Brogan, W.L., Lee, G.K.F., Sage, A.P., Kuo, B.C., Phillips, C.L., Harbor, R.D., Jacquot, R.G., McInroy, J.E., Atherton, D.P., Bay, J.S., Baumann, W.T., Chow, M-Y.  "Control Systems"
*The Electrical Engineering Handbook*
Ed. Richard C. Dorf
Boca Raton: CRC Press LLC, 2000

# 100
# Control Systems

**William L. Brogan**
*University of Nevada, Las Vegas*

**Gordon K. F. Lee**
*North Carolina State University*

**Andrew P. Sage**
*George Mason University*

**Benjamin C. Kuo**
*University of Illinois (Urbana-Champaign)*

**Charles L. Phillips**
*Auburn University*

**Royce D. Harbor**
*University of West Florida*

**Raymond G. Jacquot**
*University of Wyoming*

**John E. McInroy**
*University of Wyoming*

**Derek P. Atherton**
*University of Sussex*

**John S. Bay**
*Virginia Polytechnic Institute and State University*

**William T. Baumann**
*Virginia Polytechnic Institute and State University*

**Mo-Yuen Chow**
*North Carolina State University*

## 100.1 Models

*William L. Brogan*

A naive trial-and-error approach to the design of a control system might consist of constructing a controller, installing it into the system to be controlled, performing tests, and then modifying the controller until satisfactory performance is achieved. This approach could be dangerous and uneconomical, if not impossible. A more rational approach to control system design uses mathematical models. A *model* is a mathematical description of system behavior, as influenced by input variables or initial conditions. The model is a stand-in for the actual system during the control system design stage. It is used to predict performance; to carry out stability, sensitivity, and trade-off

# Control Mechanism
# for Rocket Apparatus

*Robert H. Goddard*
*Patented April 2, 1946*
*#2,397,657*

An excerpt from Robert Goddard's patent application:

*This invention relates to rockets and rocket craft which are propelled by combustion apparatus using liquid fuel and a liquid to support combustion, such as liquid oxygen. Such combustion apparatus is disclosed in my prior application Serial No. 327,257 filed April 1, 1940.*

*It is the general object of my present invention to provide control mechanism by which the necessary operative steps and adjustments for such mechanism will be affected automatically and in predetermined and orderly sequence.*

*To the attainment of this object, I provide control mechanism which will automatically discontinue flight in a safe and orderly manner.*

Dr. Goddard was instrumental in developing rocket propulsion in this country, both solid-fuel rocket engines and later liquid-fuel rocket motors used in missile and spaceflight applications. Goddard died in 1945, before this pivotal patent (filed June 23, 1941) on automatic control of liquid-fuel rockets was granted. He assigned half the rights to the Guggenheim Foundation in New York. (Copyright © 1995, Dewray Products, Inc. Used with permission.)
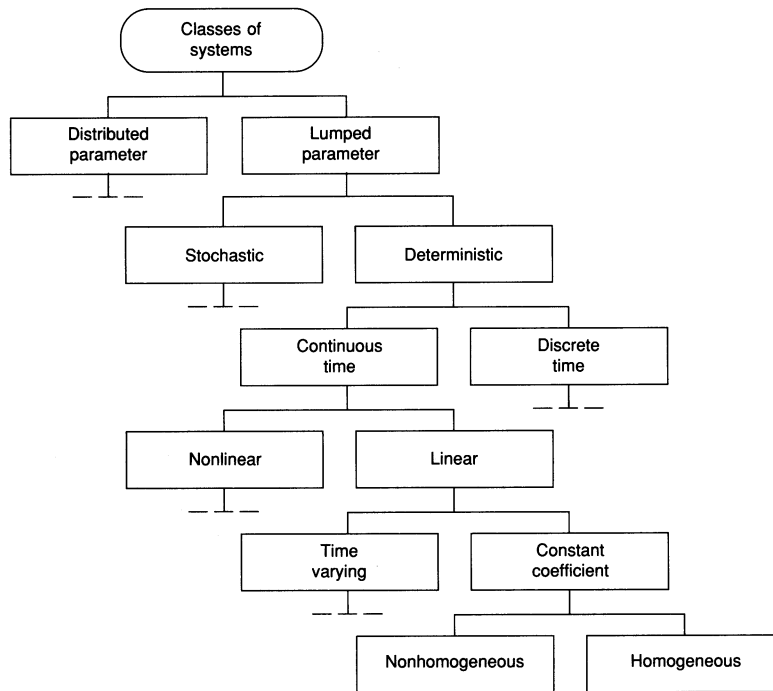
**FIGURE 100.1**  Major classes of system equations. (*Source:* W.L. Brogan, *Modern Control Theory,* 3rd ed., Englewood Cliffs, N.J.: Prentice-Hall, 1991, p. 13. With permission.)

studies; and answer various "what-if" questions in a safe and efficient manner. Of course, the validation of the model, and all conclusions derived from it, must ultimately be based upon test results with the physical hardware.

The final form of the mathematical model depends upon the type of physical system, the method used to develop the model, and mathematical manipulations applied to it. These issues are discussed next.

## Classes of Systems to Be Modeled

Most control problems are multidisciplinary. The system may consist of electrical, mechanical, thermal, optical, fluidic, or other physical components, as well as economic, biological, or ecological systems. Analogies exist between these various disciplines, based upon the similarity of the equations that describe the phenomena. The discussion of models in this section will be given in mathematical terms and therefore will apply to several disciplines.

Figure 100.1 [Brogan, 1991] shows the classes of systems that might be encountered in control systems modeling. Several branches of this tree diagram are terminated with a dashed line indicating that additional branches have been omitted, similar to those at the same level on other paths.

*Distributed parameter* systems have variables that are functions of both space and time (such as the voltage along a transmission line or the deflection of a point on an elastic structure). They are described by partial differential equations. These are often approximately modeled as a set of *lumped parameter* systems (described by ordinary differential or difference equations) by using modal expansions, finite element methods, or other approximations [Brogan, 1968]. The lumped parameter continuous-time and discrete-time families are stressed here.

## Two Major Approaches to Modeling

In principle, models of a given physical system can be developed by two distinct approaches. Figure 100.2 shows the steps involved in *analytical modeling*. The real-world system is represented by an interconnection of idealized elements. Table 100.1 [Dorf, 1989] shows model elements from several disciplines and their elemental equations. An electrical circuit diagram is a typical result of this physical modeling step (box 3 of Fig. 100.2). Application of the appropriate physical laws (Kirchhoff, Newton, etc.) to the idealized physical model (consisting of point masses, ideal springs, lumped resistors, etc.) leads to a set of mathematical equations. For a circuit these will
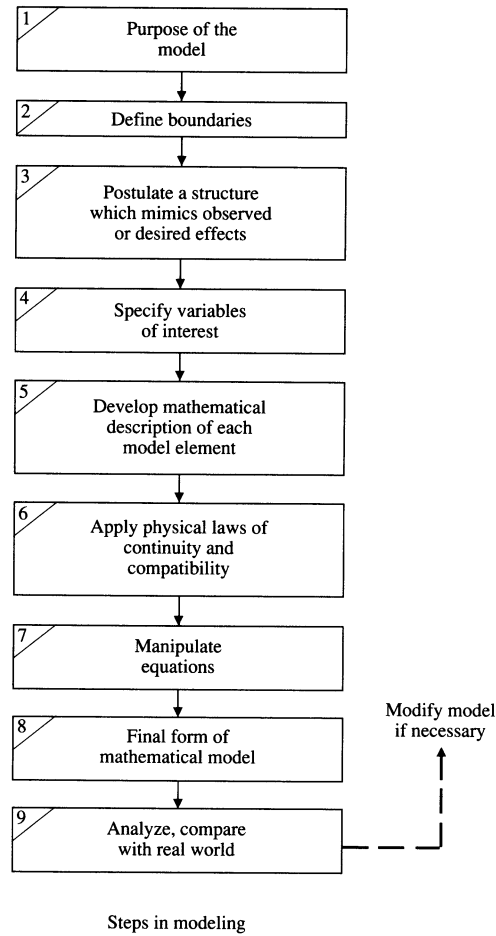
```
┌────────────────────────┐
│ 1  Purpose of the      │
│      model             │
└────────────────────────┘
            │
┌────────────────────────┐
│ 2  Define boundaries   │
└────────────────────────┘
            │
┌────────────────────────┐
│ 3  Postulate a structure│
│    which mimics observed│
│    or desired effects  │
└────────────────────────┘
            │
┌────────────────────────┐
│ 4  Specify variables   │
│      of interest       │
└────────────────────────┘
            │
┌────────────────────────┐
│ 5  Develop mathematical │
│    description of each │
│    model element       │
└────────────────────────┘
            │
┌────────────────────────┐
│ 6  Apply physical laws of│
│    continuity and      │
│    compatibility       │
└────────────────────────┘
            │
┌────────────────────────┐
│ 7  Manipulate          │
│      equations         │
└────────────────────────┘
            │
┌────────────────────────┐          Modify model
│ 8  Final form of       │          if necessary
│    mathematical model  │              ↑
└────────────────────────┘              │
            │                           │
┌────────────────────────┐              │
│ 9  Analyze, compare    │ ─ ─ ─ ─ ─ ─ ┘
│    with real world     │
└────────────────────────┘

           Steps in modeling
```

**FIGURE 100.2**  Modeling considerations. (*Source:* W.L. Brogan, *Modern Control Theory,* 3rd ed., Englewood Cliffs, N.J.: Prentice-Hall, 1991, p. 5. With permission.)

be mesh or node equations in terms of elemental currents and voltages. Box 6 of Fig. 100.2 suggests a generalization to other disciplines, in terms of continuity and compatibility laws, using through variables (generalization of current that flows through an element) and across variables (generalization of voltage, which has a differential value across an element) [Shearer et al., 1967; Dorf, 1989].
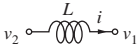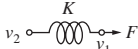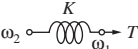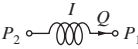
*Experimental* or *empirical* modeling typically assumes an *a priori* form for the model equations and then uses available measurements to estimate the coefficient values that cause the assumed form to best fit the data. The assumed form could be based upon physical knowledge or it could be just a credible assumption. Time-series models include autoregressive (AR) models, moving average (MA) models, and the combination, called ARMA models. All are difference equations relating the input variables to the output variables at the discrete measurement times, of the form

$$
\begin{aligned}
y(k+1) = {} & a_0 y(k) + a_1 y(k-1) + a_2 y(k-2) + \dots + a_n y(k-n) \\
& + b_0 u(k+1) + b_1 u(k) + \dots + b_p u(k+1-p) + v(k)
\end{aligned} \tag{100.1}
$$

where $v(k)$ is a random noise term. The $z$-transform transfer function relating $u$ to $y$ is

$$
\frac{y(z)}{u(z)} = \frac{b_0 + b_1 z^{-1} + \dots + b_p z^{-p}}{1 - (a_0 z^{-1} + \dots + a_{n-1} z^{-n})} = H(z) \tag{100.2}
$$

**TABLE 100.1** Summary of Describing Differential Equations for Ideal Elements

| Type of Element | Physical Element | Describing Equation | Energy $E$ or Power $\mathcal{P}$ | Symbol |
|---|---|---|---|---|
| Inductive storage | Electrical inductance | $v_{21} = L \dfrac{di}{dt}$ | $E = \dfrac{1}{2} Li^2$ | $v_2 \,\text{—}\!\!\text{\raisebox{0pt}{$\bigcirc\!\bigcirc\!\bigcirc$}}\xrightarrow{i} v_1$ |
| | Translational spring | $v_{21} = \dfrac{1}{K} \dfrac{dF}{dt}$ | $E = \dfrac{1}{2} \dfrac{F^2}{K}$ | |
| | Rotational spring | $\omega_{21} = \dfrac{1}{K} \dfrac{dT}{dt}$ | $E = \dfrac{1}{2} \dfrac{T^2}{K}$ | |
| | Fluid inertia | $P_{21} = I \dfrac{dQ}{dt}$ | $E = \dfrac{1}{2} IQ^2$ | |
| Capacitive storage | Electrical capacitance | $i = C \dfrac{dv_{21}}{dt}$ | $E = \dfrac{1}{2} Cv_{21}^2$ | |
| | Translational mass | $F = M \dfrac{dv_2}{dt}$ | $E = \dfrac{1}{2} Mv_2^2$ | |
| | Rotational mass | $T = J \dfrac{d\omega_2}{dt}$ | $E = \dfrac{1}{2} J\omega_2^2$ | |
| | Fluid capacitance | $Q = C_f \dfrac{dP_{21}}{dt}$ | $E = \dfrac{1}{2} C_f P_{21}^2$ | |
| | Thermal capacitance | $q = C_t \dfrac{d\tau_2}{dt}$ | $E = C_t \tau_2$ | |
| Energy dissipators | Electrical resistance | $i = \dfrac{1}{R} v_{21}$ | $\mathcal{P} = \dfrac{1}{R} v_{21}^2$ | |
| | Translational damper | $F = f v_{21}$ | $\mathcal{P} = f v_{21}^2$ | |
| | Rotational damper | $T = f \omega_{21}$ | $\mathcal{P} = f \omega_{21}^2$ | |
| | Fluid resistance | $Q = \dfrac{1}{R_f} P_{21}$ | $\mathcal{P} = \dfrac{1}{R_f} P_{21}^2$ | |
| | Thermal resistance | $q = \dfrac{1}{R_t} \tau_{21}$ | $\mathcal{P} = \dfrac{1}{R_t} \tau_{21}$ | |

In the MA model all $a_i = 0$. This is alternatively called an all-zero model or a finite impulse response (FIR) model. In the AR model all $b_j$ terms are zero except $b_0$. This is called an all-pole model or an infinite impulse response (IIR) model. The ARMA model has both poles and zeros and also is an IIR model [Makhoul, 1975].

Adaptive and learning control systems have an experimental modeling aspect. The data fitting is carried out on-line, in real time, as part of the system operation. The modeling described above is normally done off-line [Astrom and Wittenmark, 1989].

## Forms of the Model

Regardless of whether a model is developed from knowledge of the physics of the process or from empirical data fitting, it can be further manipulated into several different but equivalent forms. This manipulation is box 7 in Fig. 100.2. The class that is most widely used in control studies is the deterministic lumped-parameter continuous-time constant-coefficient system. A simple example has one input $u$ and one output $y$. This might be a circuit composed of one ideal source and an interconnection of ideal resistors, capacitors, and inductors.

The equations for this system might consist of a set of mesh or node equations. These could be reduced to a single $n$th-order linear ordinary differential equation by eliminating extraneous variables.

$$\frac{d^n y}{dt^n} + a_{n-1}\frac{d^{n-1}y}{dt^{n-1}} + \cdots + a_1\frac{dy}{dt} + a_0 y = b_0 u + b_1\frac{du}{dt} + \cdots + b_m\frac{d^m u}{dt^m} \qquad (100.3)$$

This $n$th-order equation can be replaced by an input-output transfer function

$$\frac{Y(s)}{U(s)} = H(s) = \frac{b_m s^m + b_{m-1}s^{m-1} + \cdots + b_1 s + b_0}{s^n + a_{n-1}s^{n-1} + \cdots + a_1 s + a_0} \qquad (100.4)$$

The inverse Laplace transform $\mathcal{L}^{-1}\{H(s)\} = h(t)$ is the system impulse response function. Alternatively, by selecting a set of $n$ internal **state variables,** Eq.(100.3) can be written as a coupled set of first-order differential equations plus an algebraic equation relating the states to the original output $y$. These equations are called state equations, and one possible choice for this example is, assuming $m = n,$

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} -a_{n-1} & 1 & 0 & 0 & \cdots & 0 \\ -a_{n-2} & 0 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ -a_1 & 0 & 0 & 0 & \cdots & 1 \\ -a_0 & 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} b_{n-1} - a_{n-1}b_n \\ b_{n-2} - a_{n-2}b_n \\ \vdots \\ b_1 - a_1 b_n \\ b_0 - a_0 b_n \end{bmatrix} u(t)$$

and

$$y(t) = [1 \ \ 0 \ \ 0 \ \ \ldots \ \ 0]\mathbf{x}(t) + b_n u(t) \qquad (100.5)$$

In matrix notation these are written more succinctly as

$$\dot{\mathbf{x}} = A\mathbf{x} + Bu \quad \text{and} \quad y = C\mathbf{x} + Du \qquad (100.6)$$

Any one of these six possible model forms, or others, might constitute the result of box 8 in Fig. 100.2. Discrete-time system models have similar choices of form, including an $n$th-order difference equation as given in Eq. (100.1) or a $z$-transform input-output transfer function as given in Eq. (100.2). A set of $n$ first-order difference equations (state equations) analogous to Eq. (100.5) or (100.6) also can be written.

Extensions to systems with $r$ inputs and $m$ outputs lead to a set of $m$ coupled equations similar to Eq. (100.3), one for each output $y_i$. These higher-order equations can be reduced to $n$ first-order state differential equations and $m$ algebraic output equations as in Eq. (100.5) or (100.6). The **A** matrix is again of dimension $n \times n,$ but **B** is now $n \times r,$ **C** is $m \times n,$ and **D** is $m \times r.$ In all previous discussions, the number of state variables, $n,$ is the order of the model. In transfer function form, an $m \times r$ matrix $H(s)$ of transfer functions will describe the input-output behavior

$$Y(s) = H(s)U(s) \qquad (100.7)$$

Other transfer function forms are also applicable, including the left and right forms of the matrix fraction description (MFD) of the transfer functions [Kailath, 1980]

$$H(s) = P(s)^{-1}N(s) \quad \text{or} \quad H(s) = N(s)P(s)^{-1} \tag{100.8}$$

Both **P** and **N** are matrices whose elements are polynomials in $s$. Very similar model forms apply to continuous-time and discrete-time systems, with the major difference being whether Laplace transform or $z$-transform transfer functions are involved.

When time-variable systems are encountered, the option of using high-order differential or difference equations versus sets of first-order state equations is still open. The system coefficients $a_i(t)$, $b_j(t)$ and/or the matrices $\mathbf{A}(t)$, $\mathbf{B}(t)$, $\mathbf{C}(t)$, and $\mathbf{D}(t)$ will now be time-varying. Transfer function approaches lose most of their utility in time-varying cases and are seldom used. With nonlinear systems all the options relating to the order and number of differential or difference equation still apply.

The form of the nonlinear state equations is

$$\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}, \mathbf{t})$$

$$y = h(\mathbf{x}, \mathbf{u}, \mathbf{t}) \tag{100.9}$$

where the nonlinear vector-valued functions $f(\mathbf{x}, \mathbf{u}, \mathbf{t})$ and $h(\mathbf{x}, \mathbf{u}, \mathbf{t})$ replace the right-hand sides of Eq. (100.6). The transfer function forms are of no value in nonlinear cases.

Stochastic systems [Maybeck, 1979] are modeled in similar forms, except the coefficients of the model and/or the inputs are described in probabilistic terms.

## Nonuniqueness

There is not a unique correct model of a given system for several reasons. The selection of idealized elements to represent the system requires judgment based upon the intended purpose. For example, a satellite might be modeled as a point mass in a study of its gross motion through space. A detailed flexible structure model might be required if the goal is to control vibration of a crucial on-board sensor. In empirical modeling, the assumed starting form, Eq. (100.1), can vary.

There is a trade-off between the complexity of the model form and the fidelity with which it will match the data set. For example, a $p$th-degree polynomial can exactly fit to $p + 1$ data points, but a straight line might be a better model of the underlying physics. Deviations from the line might be caused by extraneous measurement noise. Issues such as these are addressed in Astrom [1980].

The preceding paragraph addresses nonuniqueness in determining an input-output system description. In addition, state models developed from input-output descriptions are not unique. Suppose the transfer function of a single-input, single-output linear system is known exactly. The state variable model of this system is not unique for at least two reasons. An arbitrarily high-order state variable model can be found that will have this same transfer function. There is, however, a unique minimal or irreducible order $n_{min}$ from among all state models that have the specified transfer function. A state model of this order will have the desirable properties of **controllability** and **observability**. It is interesting to point out that the minimal order may be less than the actual order of the physical system.

The second aspect of the nonuniqueness issue relates not to order, i.e., the *number* of state variables, but to *choice* of internal variables (state variables). Mathematical and physical methods of selecting state variables are available [Brogan, 1991]. An infinite number of choices exist, and each leads to a different set {*A, B, C, D*}, called a realization. Some state variable model forms are more convenient for revealing key system properties such as stability, controllability, observability, **stabilizability,** and **detectability.** Common forms include the controllable canonical form, the observable canonical form, the Jordan canonical form, and the Kalman canonical form.

The reverse process is unique in that every valid realization leads to the same model transfer function
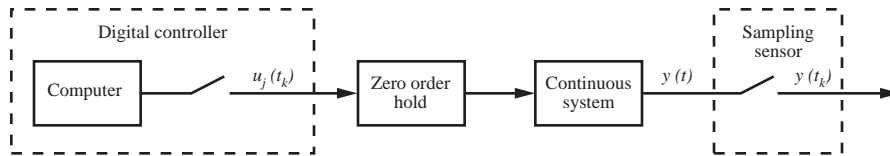
$$H(s) = C\{sI - A\}^{-1}B + D \tag{100.10}$$

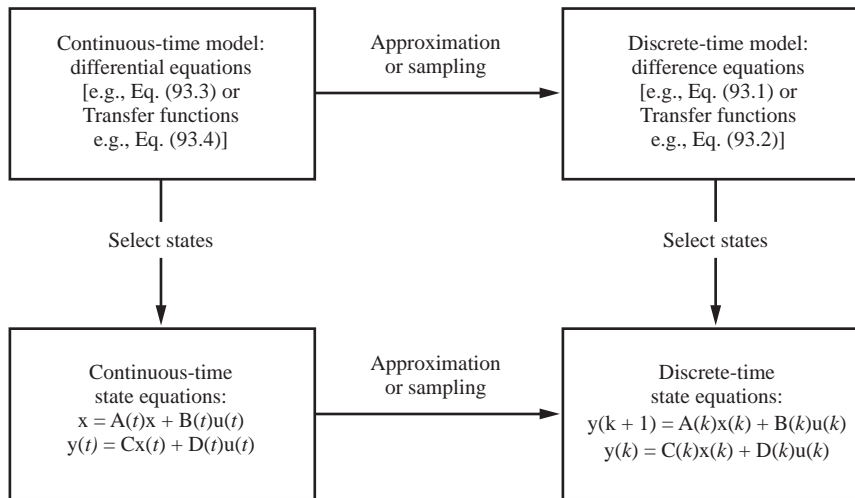**FIGURE 100.3** Digital output provided by modern sensor.



**FIGURE 100.4** State variable modeling paradigm.

## Approximation of Continuous Systems by Discrete Models

Modern control systems often are implemented digitally, and many modern sensors provide digital output, as shown in Fig. 100.3. In designing or analyzing such systems discrete-time approximate models of continuous-time systems are frequently needed. There are several general ways of proceeding, as shown in Fig. 100.4. Many choices exist for each path on the figure. Alternative choices of states or of approximation methods, such as forward or backward differences, lead to an infinite number of valid models.

## Defining Terms

**Controllability:** A property that in the linear system case depends upon the **A,B** matrix pair which ensures the existence of some control input that will drive any arbitrary initial state to zero in finite time.

**Detectability:** A system is detectable if all its unstable modes are observable.

**Observability:** A property that in the linear system case depends upon the **A,C** matrix pair which ensures the ability to determine the initial values of all states by observing the system outputs for some finite time interval.

**Stabilizable:** A system is stabilizable if all its unstable modes are controllable.

**State variables:** A set of variables that completely summarize the system's status in the following sense. If all states $x_i$ are known at time $t_0$, then the values of all states and outputs can be determined uniquely for any time $t_1 > t_0$, provided the inputs are known from $t_0$ onward. State variables are components in the state vector. State space is a vector space containing the state vectors.

## Related Topic

6.1 Definitions and Properties

## References

K.J. Astrom, "Maximum likelihood and prediction error methods," *Automatica,* vol. 16, pp. 551–574, 1980.

K.J. Astrom and B. Wittenmark, *Adaptive Control,* Reading, Mass.: Addison-Wesley, 1989.

W.L. Brogan, "Optimal control theory applied to systems described by partial differential equations," in *Advances in Control Systems,* vol. 6, C. T. Leondes (ed.), New York: Academic Press, 1968, chap. 4.

W.L. Brogan, *Modern Control Theory,* 3rd ed., Englewood Cliffs, N.J.: Prentice-Hall, 1991.

R.C. Dorf, *Modern Control Systems,* 5th ed., Reading, Mass.: Addison-Wesley, 1989.

T. Kailath, *Linear Systems,* Englewood Cliffs, N.J.: Prentice-Hall, 1980.

J. Makhoul, "Linear prediction: A tutorial review," *Proc. IEEE,* vol. 63, no. 4, pp. 561–580, 1975.

P.S. Maybeck, *Stochastic Models, Estimation and Control,* vol. 1, New York: Academic Press, 1979.

J.L. Shearer, A.T. Murphy, and H.H. Richardson, *Introduction to Dynamic Systems,* Reading, Mass.: Addison-Wesley, 1967.

## Further Information

The monthly *IEEE Control Systems Magazine* frequently contains application articles involving models of interesting physical systems.

The monthly *IEEE Transactions on Automatic Control* is concerned with theoretical aspects of systems. Models as discussed here are often the starting point for these investigations.

*Automatica* is the source of many related articles. In particular an extended survey on system identification is given by Astrom and Eykhoff in vol. 7, pp. 123–162, 1971.

Early developments of the state variable approach are given by R. E. Kalman in "Mathematical description of linear dynamical systems," *SIAM J. Control Ser.,* vol. A1, no. 2, pp. 152–192, 1963.

# 100.2  Dynamic Response

*Gordon K. F. Lee*

## Computing the Dynamic System Response

Consider a linear time-invariant dynamic system represented by a differential equation form

$$
\frac{d^n y(t)}{dt^n} + a_{n-1}\frac{d^{n-1}y(t)}{dt^{n-1}} + \cdots + a_1\frac{dy(t)}{dt} + a_0 y(t)
$$
$$
= b_m\frac{d^m f(t)}{dt^m} + \cdots + b_1\frac{df(t)}{dt} + b_0 f(t)
$$

(100.11)

where $y(t)$ and $f(t)$ represent the output and input, respectively, of the system.

Let $p^k(\cdot) \overset{\Delta}{=} (d^k/dt^k)(\cdot)$ define the differential operator so that (100.11) becomes

$$
(p^n + a_{n-1}p^{n-1} + \ldots + a_1 p + a_0)y(t) = (b_m p^m + \ldots + b_1 p + b_0)f(t)
$$

(100.12)

The solution to (100.11) is given by

$$
y(t) = y_S(t) + y_I(t)
$$

(100.13)

where $y_S(t)$ is the **zero-input response,** or that part of the response due to the initial conditions (or states) only, and $y_I(t)$ is the **zero-state response,** or that part of the response due to the input $f(t)$ only.

## Zero-Input Response: $y_S(t)$

Here $f(t) = 0$, and thus (100.11) becomes

$$(p^n + a_{n-1}p^{n-1} + \dots + a_1p + a_0)y(t) = 0 \tag{100.14}$$

That is,

$$D(p)y(t) = 0$$

The roots of $D(p) = 0$ can be categorized as either distinct or multiple. That is, in general,

$$D(p) = \prod_{i=1}^{q}(p - \lambda_i)^{k_i}\prod_{i=1}^{r}(p - \lambda_{q+i})$$

where there are $r$ distinct roots and $q$ sets of multiple roots (each set has multiplicity $k_i$). Note that $r + \sigma = n$, where $\sigma \overset{\Delta}{=} \sum_{i=1}^{q}k_i$. Each distinct root contributes a term to $y_S(t)$ of the form $c_i e^{\lambda_i t}$, where $c_i$ is a constant, while each set of multiple roots contributes a set of terms to $y_S(t)$ of the form $\sum_{j=0}^{k_i-1}c_{i,j}t^j e^{\lambda_i t}$, where $c_{i,j}$ is some constant. Thus, the zero-input response is given by

$$y_S(t) = \sum_{i=1}^{q}\sum_{j=0}^{k_i-1}c_{i,j}t^j e^{\lambda_i t} + \sum_{i=1}^{r}c_{\sigma+i}e^{\lambda_{\sigma+i}t} \tag{100.15}$$

The coefficients $c_{i,j}$ and $c_{\sigma+i}$ are selected to satisfy the initial conditions.

## Special Case

If all the roots of $D(p) = 0$ are distinct and the initial conditions for (100.11) are given by

$$\left\{y(0), \frac{dy(0)}{dt}, \dots, \frac{d^{n-1}y(0)}{dt^{n-1}}\right\}$$

then the coefficients of (100.15) are given by the solution of

$$\begin{bmatrix} y(0) \\ \dfrac{dy(0)}{dt} \\ \vdots \\ \dfrac{d^{n-1}y(0)}{dt^{n-1}} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \lambda_1 & \lambda_2 & \cdots & \lambda_n \\ \vdots & \vdots & \vdots & \vdots \\ \lambda_1^{(n-1)} & \lambda_2^{(n-1)} & \cdots & \lambda_n^{(n-1)} \end{bmatrix}\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} \tag{100.16}$$

## Zero-State Response: $y_I(t)$

Here the initial conditions are made identically zero. Observing (100.11), let

$$H(p) = \frac{b_m p^m + \cdots + b_1 p + b_0}{p^n + a_{n-1}p^{n-1} + \cdots + a_1 p + a_0}$$

denote a rational function in the $p$ operator. Consider using partial-fraction expansion on $H(p)$ as

$$H(p) = \sum_{i=1}^{q} \sum_{j=1}^{k_i} \frac{g_{i,j}}{(p - \lambda_i)^j} + \sum_{i=1}^{r} \frac{g_{\sigma+i}}{p - \lambda_{q+i}} \qquad (100.17)$$

when the first term corresponds to the sets of multiple roots and the second term corresponds to the distinct roots.

Note the constant residuals are computed as

$$g_{\sigma+i} = \left[ (p - \lambda_{q+i}) H(p) \right]_{p = \lambda_{q+i}}$$

and

$$g_{i,j} = \frac{1}{(k_i - j)!} \frac{d^{(k_i - j)}}{dp^{(k_i - j)}} \left\{ (p - \lambda_i)^{k_i} H(p) \right\} \bigg|_{p = \lambda_i}$$

Then

$$h(t) = \sum_{i=1}^{q} \sum_{j=1}^{k_i} \frac{g_{i,j}}{(j-1)!} t^{j-1} e^{\lambda_i t} + \sum_{i=1}^{r} g_{\sigma+i} e^{\lambda_{\sigma+i} t} \qquad (100.18)$$

is the **impulse response** of the system (100.11). Then the zero-state response is given by

$$y_I(t) = \int_0^t f(\tau) h(t - \tau) \, d\tau \qquad (100.19)$$

that is, $y_I(t)$ is the time convolution between input $f(t)$ and impulse response $h(t)$. In some instances, it may be easier to find $y_s(t)$ and $y_I(t)$ using Laplace Transform methods.

## Measures of the Dynamic System Response

Several measures may be employed to investigate dynamic response performance. These include:

1. Speed of the response—how quickly does the system reach its final value
2. Accuracy—how close is the final response to the desired response
3. Relative stability—how stable is the system or how close is the system to instability
4. Sensitivity—what happens to the system response if the system parameters change

Objectives 3 and 4 may be analyzed by frequency domain methods (Section 100.3). Time-domain measures classically analyze the dynamic response by partitioning the total response into its steady-state (objective 2) and transient (objective 1) components. The **steady-state response** is that part of the response which remains as time approaches infinity; the **transient response** is that part of the response which vanishes as time approaches infinity.

### Measures of the Steady-State Response

In the steady state, the accuracy of the time response is an indication of how well the dynamic response follows a desired time trajectory. Usually a test signal (reference signal) is selected to measure accuracy. Consider Fig. 100.5. In this configuration, the objective is to force $y(t)$ to track a reference signal $r(t)$ as close as possible.
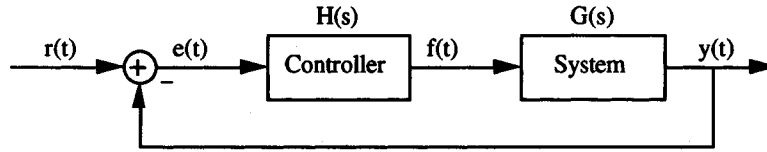
**FIGURE 100.5** A tracking controller configuration.

**TABLE 100.2** Steady-State Error Constants

| $r(t)$<br>Test Signal | $eSS(t)$ | Error<br>Constant |
|---|---|---|
| $Ru(t)$: step function | $\dfrac{R}{1 + K_p}$ | $K_p = \lim\limits_{s \to 0} G(s)H(s)$ |
| $Rtu(t)$: ramp function | $\dfrac{R}{K_v}$ | $K_v = \lim\limits_{s \to 0} sG(s)H(s)$ |
| $\dfrac{R}{2}t^2u(t)$: parabolic function | $\dfrac{R}{K_a}$ | $K_a = \lim\limits_{s \to 0} s^2G(s)H(s)$ |

The **steady-state error** is a measure of the accuracy of the output $y(t)$ in tracking the reference input $r(t)$. Other configurations with different performance measures would result in other definitions of the steady-state error between two signals.

From Fig. 100.5, the error $e(t)$ is

$$e(t) = r(t) - y(t) \tag{100.20}$$

and the steady-state error is

$$e_{SS}(t) = \lim_{t \to \infty} e(t) = \lim_{s \to \infty} sE(s) \tag{100.21}$$

assuming the limits exists, where $E(s)$ is the Laplace transform of $e(t)$, and $s$ is the Laplacian operator. With $G(s)$ the transfer function of the system and $H(s)$ the transfer function of the controller, the transfer function between $y(t)$ and $r(t)$ is found to be

$$T(s) = \frac{G(s)H(s)}{1 + G(s)H(s)} \tag{100.22}$$

with

$$E(s) = \frac{R(s)}{1 + G(s)H(s)} \tag{100.23}$$

Direct application of the steady-state error for various inputs yields Table 100.2. Note $u(t)$ is the unit step function. This table can be extended to an $m$th-order input in a straightforward manner. Note that for $e_{SS}(t)$ to go to zero with a reference signal $Ct^mu(t)$, the term $G(s)H(s)$ must have at least $m$ poles at the origin (a type $m$ system).
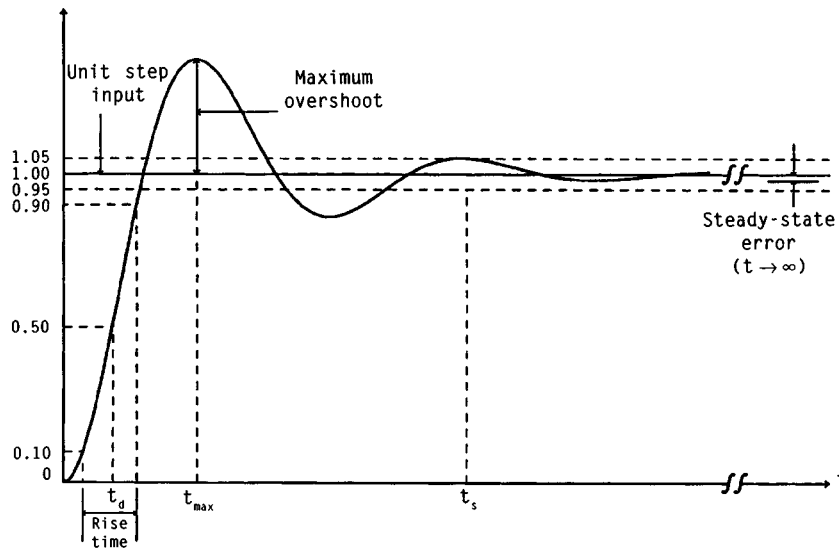
**FIGURE 100.6** Step response.

### Measures of the Transient Response

In general, analysis of the transient response of a dynamic system to a reference input is difficult. Hence formulating a standard measure of performance becomes complicated. In many cases, the response is dominated by a pair of poles and thus acts like a second-order system.

Consider a reference unit step input to a dynamic system (Fig. 100.6). Critical parameters that measure transient response include:

1. *M:* maximum overshoot
2. % overshoot = $M/A \times 100\%$, where $A$ is the final value of the time response
3. $t_d$: delay time—the time required to reach 50% of $A$
4. $t_r$: rise time—the time required to go from 10% of $A$ to 90% of $A$
5. $t_s$: settling time—the time required for the response to reach and stay within 5% of $A$

To calculate these measures, consider a second-order system

$$T(s) = \frac{\omega_n^2}{s^2 + 2\xi\omega_n s + \omega_n^2} \tag{100.24}$$

where $\xi$ is the damping coefficient and $\omega_n$ is the natural frequency of oscillation.

For the range $0 < \xi < 1$, the system response is *underdamped,* resulting in a damped oscillatory output. For a unit step input, the response is given by

$$y(t) = 1 + \frac{e^{-\xi\omega_n t}}{\sqrt{1 - \xi^2}} \sin\left(\omega_n\sqrt{1 - \xi^2}\, t - \tan^{-1}\frac{\sqrt{1 - \xi^2}}{-\xi}\right) \quad (0 < \xi < 1) \tag{100.25}$$

The eigenvalues (poles) of the system [roots of the denominator of $T(s)$] provide some measure of the time constants of the system. For the system under study, the eigenvalues are at

$$-\xi\omega_n \pm j\omega_n\sqrt{1 - \xi^2} \qquad \text{where} \qquad j \triangleq \sqrt{-1}$$

**FIGURE 100.7**  Effect of the damping coefficient on the dynamic response.



**FIGURE 100.8**  Effect of the natural frequency of oscillation on the dynamic response.

From the expression of $y(t)$, one sees that the term $\xi\omega_n$ affects the rise time and exponential decay time. The effects of the damping coefficient on the transient response are seen in Fig. 100.7.

The effects of the natural frequency of oscillation $\omega_n$ of the transient response can be seen in Fig. 100.8. As $\omega_n$ increases, the frequency of oscillation increases.

For the case when $0 < \xi < 1$, the underdamped case, one can analyze the critical transient response parameters.

To measure the peaks of Fig. 100.6, one finds

$$y_{peak}(t) = 1 + (-1)^{n-1} \exp \frac{-n\pi\xi}{\sqrt{1-\xi^2}} \qquad n = 0, 1, \ldots \qquad (100.26)$$

occurring at

$$t = \frac{n\pi}{\omega_n\sqrt{1-\xi^2}} \qquad \begin{array}{l} n: \text{odd (overshoot)} \\ n: \text{even (undershoot)} \end{array} \qquad (100.27)$$

Hence

$$y_{max} = 1 + \exp \frac{-\pi\xi}{\sqrt{1-\xi^2}} \qquad (100.28)$$

occurring at

$$t_{max} = \frac{\pi}{\omega_n\sqrt{1-\xi^2}} \qquad (100.29)$$

With these parameters, one finds

$$t_d \approx \frac{1 + 0.7\xi}{\omega_n}$$

$$t_r \approx \frac{1 + 1.1\xi + 1.4\xi^2}{\omega_n}$$

and

$$t_s \approx \frac{3}{\xi\omega_n}$$

Note that increasing $\xi$ decreases the % overshoot and decreases the settling time but increases $t_d$ and $t_r$.

When $\xi = 1$, the system has a double pole at $-\omega_n$, resulting in a *critically damped* response. This is the point when the response just changes from oscillatory to exponential in form. For a unit step input, the response is given by

$$y(t) = 1 - e^{-\omega nt}(1 + \omega_n t) \quad (\xi = 1) \qquad (100.30)$$

For the range $\xi > 1$, the system is overdamped due to two real system poles. For a unit step input, the response is given by

$$y(t) = 1 + \frac{1}{c_1 - c_2}\left(\frac{1}{c_1} e^{c_1\omega_n t} - \frac{1}{c_2} e^{c_2\omega_n t}\right) \quad (\xi > 1)$$

$$c_1 = -\xi + \sqrt{\xi^2 - 1} \qquad c_2 = -\xi - \sqrt{\xi^2 - 1} \qquad (100.31)$$

Finally, when $\xi = 0$, the response is purely sinusoidal. For a unit step, the response is given by

$$y(t) = 1 - \cos \omega_n t \quad (\xi = 0) \qquad (100.32)$$

## Defining Terms

**Impulse response:** The response of a system when the input is an impulse function.
**Steady-state error:** The difference between the desired reference signal and the actual signal in steady-state, i.e., when time approaches infinity.
**Steady-state response:** That part of the response which remains as time approaches infinity.
**Transient response:** That part of the response which vanishes as time approaches infinity.
**Zero-input response:** That part of the response due to the initial condition only.
**Zero-state response:** That part of the response due to the input only.

## Related Topics

6.1 Definitions and Properties • 7.1 Introduction • 112.2 A Brief History of CACSD

## Further Information

J.J. D'Azzo and C.H. Harpis, *Linear Control System Analysis and Design,* New York: McGraw-Hill, 1981.
R.C. Dorf, *Modern Control Systems,* 5th ed., Reading, Mass.: Addison-Wesley, 1989.
M.E. El-Hawary, *Control Systems Engineering,* Reston, Va.: Reston, 1984.
G.H. Hostetter, C. J. Savant, Jr., and R. T. Stefani, *Design of Feedback Control Systems,* Philadelphia: Saunders, 1989.
B.C. Kuo, *Automatic Control Systems,* Englewood Cliffs, N.J.: Prentice-Hall, 1987.
K. Ogata, *Modern Control Engineering,* Englewood Cliffs, N.J.: Prentice-Hall, 1970.
N.K. Sinha, *Control Systems,* New York: Holt, 1986.

## 100.3  Frequency Response Methods: Bode Diagram Approach

*Andrew P. Sage*

Our efforts in this section are concerned with analysis and design of linear control systems by frequency response methods. Design generally involves trial-and-error repetition of analysis until a set of design **specifications** has been met. Thus, analysis methods are most useful in the design process, which is one phase of the **systems engineering** life cycle [Sage, 1992]. We will discuss one design method based on **Bode diagrams.** We will discuss the use of both simple **series equalizers** and composite equalizers as well as the use of minor-loop feedback in systems design.

Figure 100.9 presents a flowchart of the frequency response method design process and indicates the key role of analysis in linear systems control design. The flowchart of Fig. 100.9 is applicable to control system design methods in general. There are several iterative loops, generally calling for trial-and-error efforts, that comprise the suggested design process. An experienced designer will often be able, primarily due to successful prior experience, to select a system structure and generic components such that the design specifications can be met with no or perhaps a very few iterations through the iterative loop involving adjustment of equalizer or compensation parameters to best meet specifications.

If the parameter optimization, or parameter refinement such as to lead to maximum phase margin, approach shows the specifications cannot be met, we are then assured that no **equalizer** of the specific form selected will meet specifications. The next design step, if needed, would consist of modification of the equalizer form or
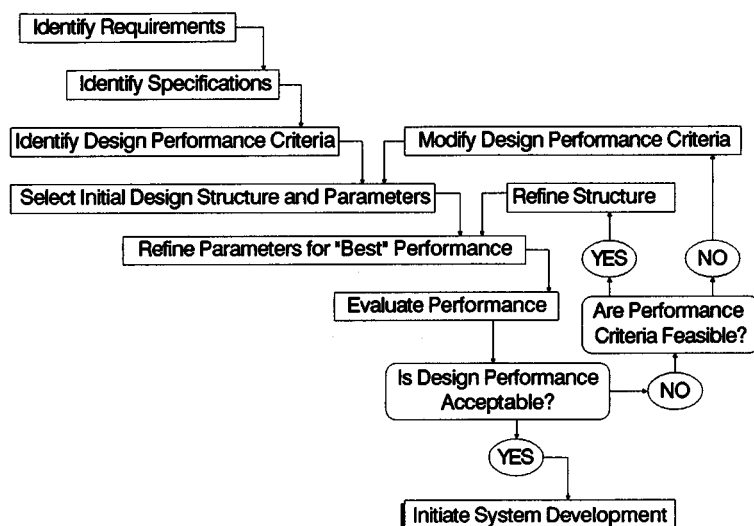
**FIGURE 100.9**  System design life cycle for frequency-response-based design.

structure and repetition of the analysis process to determine equalizer parameter values to best meet specifications. If specifications still cannot be met, we will usually next modify generic fixed components used in the system. This iterative design and analysis process is again repeated. If no reasonable fixed components can be obtained to meet specifications, then structural changes in the proposed system are next contemplated. If no structure can be found that allows satisfaction of specifications, either the client must be requested to relax the frequency response specifications or the project may be rejected as infeasible using present technology. As we might suspect, economics will play a dominant role in this design process. Changes made due to iteration in the inner loops of Fig. 100.9 normally involve little additional costs, whereas those made due to iterations in the outer loops will often involve major cost changes.

## Frequency Response Analysis Using the Bode Diagram

The steady-state response of a stable linear constant-coefficient system has particular significance, as we know from an elementary study of electrical networks and circuits and of dynamics. We consider a stable linear system with input-output transfer function

$$H(s) = \frac{Z(s)}{U(s)}$$

We assume a sinusoidal input $u(t) = \cos \omega t$ so that we have for the Laplace transform of the system output

$$Z(s) = \frac{sH(s)}{s^2 + \omega^2}$$

We expand this ratio of polynomials using the partial-fraction approach and obtain

$$Z(s) = F(s) + \frac{a_1}{s + j\omega} + \frac{a_2}{s - j\omega}$$

In this expression, $F(s)$ contains all the poles of $H(s)$. All of these lie in the left half plane since the system, represented by $H(s)$, is assumed to be stable. The coefficients $a_1$ and $a_2$ are easily determined as

$$a_1 = \frac{H(-j\omega)}{2}$$

$$a_2 = \frac{H(j\omega)}{2}$$

We can represent the complex transfer function $H(j\omega)$ in either of two forms,

$$H(j\omega) = B(\omega) + jC(\omega)$$

$$H(-j\omega) = B(\omega) - jC(\omega)$$

The inverse Laplace transform of the system transfer function will result in a transient term due to the inverse transform of $F(s)$, which will decay to zero as time progresses. A steady-state component will remain, and this is, from the inverse transform of the system equation, given by

$$z(t) = a_1 e^{-j\omega t} + a_2{}^{j\omega t}$$

We combine several of these relations and obtain the result

$$z(t) = B(\omega)\left(\frac{e^{j\omega t} + e^{-j\omega t}}{2}\right) - C(\omega)\left(\frac{e^{j\omega t} - e^{-j\omega t}}{2j}\right)$$

This result becomes, using the Euler identity,[1]

$$z(t) = B(\omega)\,\cos\omega t - C(\omega)\,\sin\omega t$$

$$= [B^2(\omega) + C^2(\omega)]^{1/2}\cos(\omega + \beta)$$

$$= |H(j\omega)|\,\cos(\omega t + \beta)$$

where $\tan \beta(\omega) = C(\omega)/B(\omega)$.

As we see from this last result, there is a very direct relationship between the transfer function of a linear constant-coefficient system, the time response of a system to any known input, and the sinusoidal steady-state response of the system. We can always determine any of these if we are given any one of them. This is a very important result. This important conclusion justifies a design procedure for linear systems that is based only on sinusoidal steady-state response, as it is possible to determine transient responses, or responses to any given system input, from a knowledge of steady-state sinusoidal responses, at least in theory. In practice, this might be rather difficult computationally without some form of automated assistance.

## Bode Diagram Design-Series Equalizers

In this subsection we consider three types of series equalization:

1. Gain adjustment, normally attenuation by a constant at all frequencies
2. Increasing the phase lead, or reducing the phase lag, at the **crossover frequency** by use of a phase **lead network**
3. Attenuation of the gain at middle and high frequencies such that the crossover frequency will be decreased to a lower value where the phase lag is less, by use of a **lag network**

---

[1]The Euler identity is $e^{j\omega t} = \cos \omega t + j\sin\omega t$.
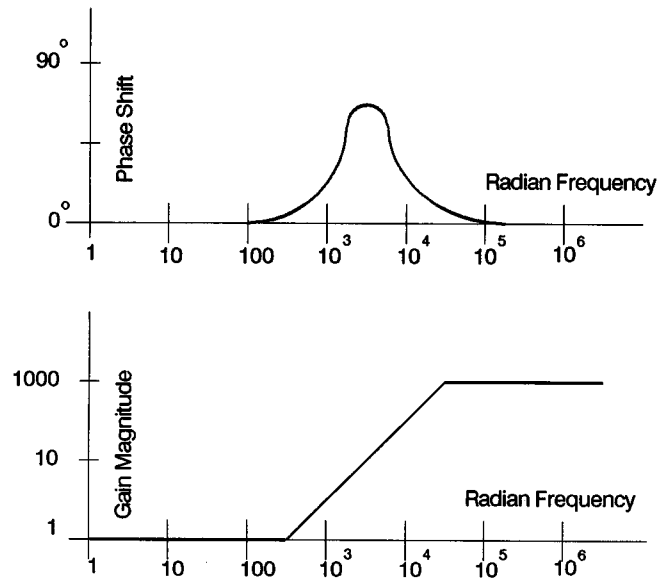
---

**FIGURE 100.10** Phase shift and gain curves for a simple lead network.

In the subsection that follows this, we will first consider use of a composite or **lag-lead network** near crossover to attenuate gain only to reduce the crossover frequency to a value where the phase lag is less. Then we will consider more complex composite equalizers and state some general guidelines for Bode diagram design. Here, we will use Bode diagram frequency domain design techniques to develop a design procedure for each of three elementary types of series equalization.

## Gain Reduction

Many linear control systems can be made sufficiently stable merely by reduction of the open-loop system gain to a sufficiently low value. This approach ignores all performance specifications, however, except that of phase margin (PM) and is, therefore, usually not a satisfactory approach. It is a very simple one, however, and serves to illustrate the approach to be taken in more complex cases.

The following steps constitute an appropriate Bode diagram design procedure for compensation by gain adjustment:
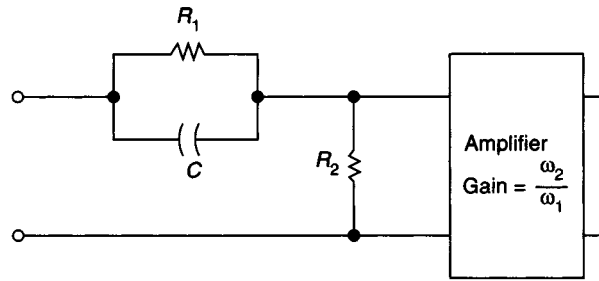
1. Determine the required PM and the corresponding phase shift $\beta_c = -\pi + PM$.
2. Determine the frequency $\omega_c$ at which the phase shift is such as to yield the phase shift at crossover required to give the desired PM.
3. Adjust the gain such that the actual crossover frequency occurs at the value computed in step 2.

## Phase-Lead Compensation

In compensation using a phase-lead network, we increase the phase lead at the crossover frequency such that we meet a performance specification concerning phase shift. A phase-lead-compensating network transfer function is

$$G_c(s) = \left(1 + \frac{s}{\omega_1}\right)\bigg/\left(1 + \frac{s}{\omega_2}\right) \qquad \omega_1 < \omega_2$$

Figure 100.10 illustrates the gain versus frequency and phase versus frequency curves for a simple lead network with the transfer function of the foregoing equation. The maximum phase lead obtainable from a phase-lead network depends upon the ratio $\omega_2/\omega_1$ that is used in designing the network. From the expression for the phase shift of the transfer function for this system, which is given by

$$G_c = \frac{1 + \dfrac{s}{\omega_1}}{1 + \dfrac{s}{\omega_2}}$$

$$\omega_1 = \frac{1}{R_1 C}, \qquad \omega_2 = \left(1 + \frac{R_1}{R_2}\right)\omega_1$$

**FIGURE 100.11**   A simple electrical lead network.

$$\beta = \tan^{-1}\frac{\omega}{\omega_1} - \tan^{-1}\frac{\omega}{\omega_2}$$

we see that the maximum amount of phase lead occurs at the point where the first derivative with respect to frequency is zero, or

$$\left.\frac{d\beta}{d\omega}\right|_{\omega=\omega_m} = 0$$

or at the frequency where

$$\omega_m = (\omega_1 \omega_2)^{0.5}$$

This frequency is easily seen to be at the center of the two break frequencies for the lead network on a Bode log asymptotic gain plot. It is interesting to note that this is exactly the same frequency that we would obtain using an arctangent approximation[2] with the assumption that $\omega_1 < \omega < \omega_2$.

There are many ways of realizing a simple phase-lead network. All methods require the use of an active element since the gain of the lead network at high frequencies is greater than 1. A simple electrical network realization is shown in Fig. 100.11.

We now consider a simple design example. Suppose that we have an open-loop system with transfer function

$$G_f(s) = \frac{10^4}{s^2}$$

It turns out that this is often called a type-two system due to the presence of the double integration. This system will have a steady-state error of zero for a constant acceleration input. The crossover frequency, that is to say

---

[2]The arctangent approximation is $\tan^{-1}(\omega/\alpha) = \omega/\alpha$ for $\omega < \alpha$ and $\tan^{-1}(\omega/\alpha) = \pi/2 - \alpha/\omega$ for $\omega > \alpha$. This approximation is rather easily obtained through use of a Taylor series approximation.
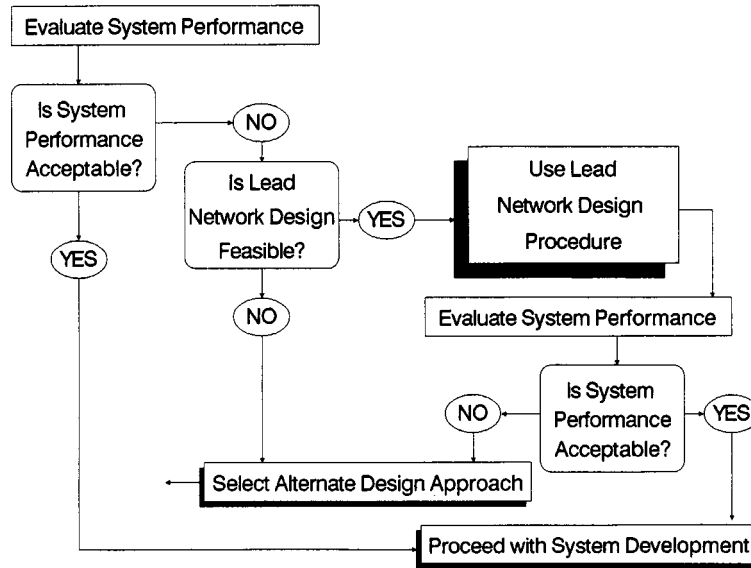
**FIGURE 100.12** Life cycle of frequency domain design incorporating lead network compensation.

the frequency where the magnitude of the open-loop gain is 1, is 100 rad/s. The PM for the system without equalization is zero. We will design a simple lead network compensation for this zero PM system. If uncompensated, the closed-loop transfer function will be such that the system is unstable and any disturbance at all will result in a sinusoidally oscillating output.

The asymptotic gain diagram for this example is easily obtained from the open-loop transfer function

$$G_f(s)G_c(s) = \frac{K(1 + s/\omega_1)}{(1 + s/\omega_2)s^2}$$

and we wish to select the break frequencies $\omega_1$ and $\omega_2$ such that the phase shift at crossover is maximum. Further, we want this maximum phase shift to be such that we obtain the specified PM. We use the procedure suggested in Fig. 100.12.

Since the crossover frequency is such that $\omega_1 < \omega_c < \omega_2$, we have for the arctangent approximation to the phase shift in the vicinity of the crossover frequency

$$\beta(\omega) = -\pi + \tan^{-1}\frac{\omega}{\omega_1} - \tan^{-1}\frac{\omega}{\omega_2}$$

$$\approx \frac{-\pi}{2} - \frac{\omega_1}{\omega} - \frac{\omega}{\omega_2}$$

In order to maximize the phase shift at crossover, we set

$$\left.\frac{d\beta}{d\omega}\right|_{\omega=\omega_m} = 0$$

and obtain as a result

$$\omega_m = (\omega_1\,\omega_2)^{0.5}$$

We see that the crossover frequency obtained is halfway between the two break frequencies $\omega_1$ and $\omega_2$ on a logarithmic frequency coordinate. The phase shift at this optimum value of crossover frequency becomes

$$\beta_c = \beta(\omega_c) = \frac{-\pi}{2} - 2\left(\frac{\omega_1}{\omega_2}\right)^{0.5}$$

For a PM of $-3\pi/4$, for example, we have $-3\pi/4 = -\pi/2 - 2(\omega_1/\omega_2)^{0.5}$, and we obtain $\omega_1/\omega_2 = 0.1542$ as the ratio of frequencies. We see that we have need for a lead network with a gain of $\omega_1/\omega_2 = 6.485$. The gain at the crossover frequency is 1, and from the asymptotic gain approximation that is valid for $\omega_1 < \omega < \omega_2$, we have the expressions $|G(j\omega)| = K/\omega\omega_1$ and $|G(j\omega_c)| = 1 = K/\omega_c\omega_1$ which for a known $K$ can be solved for $\omega_c$ and $\omega_1$.

Now that we have illustrated the design computation with a very simple example, we are in a position to state some general results. In the direct approach to design for a specified PM we assume a single lead network equalizer such that the open-loop system to transfer function results. This approach to design results in the following steps that are applicable for Bode diagram design to achieve maximum PM within an experientially determined control system structure that comprises a fixed plant and a compensation network with adjustable parameters:

1. We find an equation for the gain at the crossover frequency in terms of the compensated open-loop system break frequency.
2. We find an equation of the phase shift at crossover.
3. We find the relationship between equalizer parameters and crossover frequency such that the phase shift at crossover is the maximum possible and a minimum of additional gain is needed.
4. We determine all parameter specifications to meet the PM specifications.
5. We check to see that all design specifications have been met. If they have not, we iterate the design process.

Figure 100.12 illustrates the steps involved in implementing this frequency domain design approach.

**Phase-Lag Compensation**

In the phase-lag-compensation frequency domain design approach, we reduce the gain at low frequencies such that crossover, the frequency where the gain magnitude is 1, occurs before the phase lag has had a chance to become intolerably large. A simple single-stage phase-lag-compensating network transfer function is

$$G_c(s) = \frac{1 + s/\omega_2}{1 + s/\omega_1} \qquad \omega_1 < \omega_2$$

Figure 100.13 illustrates the gain and phase versus frequency curves for a simple lag network with this transfer function. The maximum phase lag obtainable from a phase-lag network depends upon the ratio $\omega_2/\omega_1$ that is used in designing the network. From the expression for the phase shift of this transfer function,

$$\beta = \tan^{-1}\frac{\omega}{\omega_2} - \tan^{-1}\frac{\omega}{\omega_1}$$

we see that maximum phase lag occurs at that frequency $\omega_c$ where $d\beta/d\omega = 0$. We obtain for this value

$$\omega_m = (\omega_1\omega_2)^{0.5}$$

which is at the center of the two break frequencies for the lag network when the frequency response diagram is illustrated on a Bode diagram log-log asymptotic gain plot.

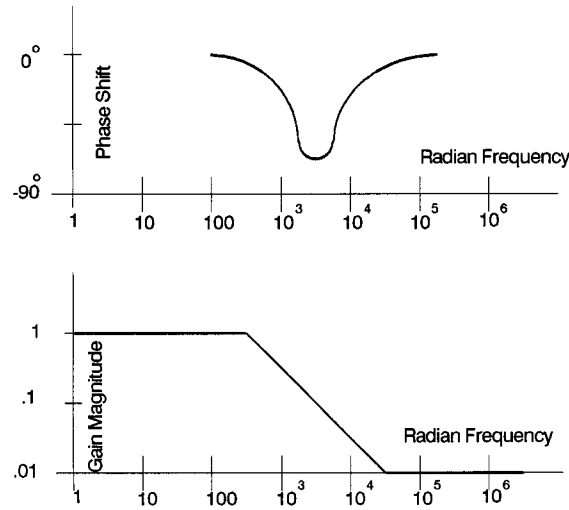The maximum value of the phase lag obtained at $\omega = \omega_m$ is

**FIGURE 100.13** Phase shift and gain curves for a simple lag network.

$$\beta_m(\omega_m) = \frac{\pi}{2} - 2\,\tan^{-1}\left(\frac{\omega_2}{\omega_1}\right)^{0.5}$$

$$= \frac{\pi}{2} - 2\,\tan^{-1}\left(\frac{\omega_1}{\omega_2}\right)^{0.5}$$

which can be approximated in a more usable form, using the arctangent approximation, as

$$\beta_m(\omega_m) \approx \frac{\pi}{2}\sqrt{\frac{\omega_2}{\omega_1}}$$

The attenuation of the lag network at the frequency of minimum phase shift, or maximum phase lag, is obtained from the asymptotic approximation as

$$\left|G_c(\omega_m)\right| = \left(\frac{\omega_1}{\omega_2}\right)^{0.5}$$

Figure 100.13 presents a curve of attenuation magnitude obtainable at the frequency of maximum phase lag and the amount of the phase lag for various ratios $\omega_2/\omega_1$ for this simple lag network.

There are many ways to physically realize a lag network transfer function. Since the network only attenuates at some frequencies, as it never has a gain greater than 1 at any frequency, it can be realized with passive components only. Figure 100.14 presents an electrical realization of the simple lag network. Figure 100.15 presents a flowchart illustrating the design procedure envisioned here for lag network design. This is conceptually very similar to that for a lead network and makes use of the five-step parameter optimization procedure suggested earlier.

The object of lag network design is to reduce the gain at frequencies lower than the original crossover frequency in order to reduce the open-loop gain to unity before the phase shift becomes so excessive that the system PM is too small. A disadvantage of lag network compensation is that the attenuation introduced reduces
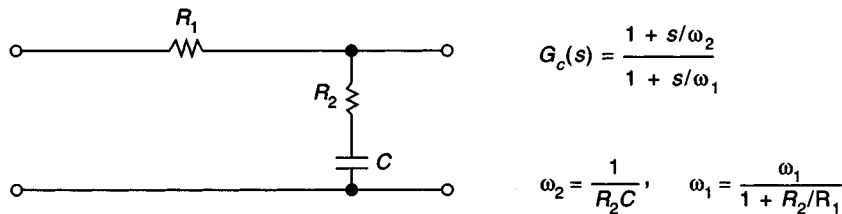
$$G_c(s) = \frac{1 + s/\omega_2}{1 + s/\omega_1}$$

$$\omega_2 = \frac{1}{R_2 C}, \qquad \omega_1 = \frac{\omega_1}{1 + R_2/R_1}$$

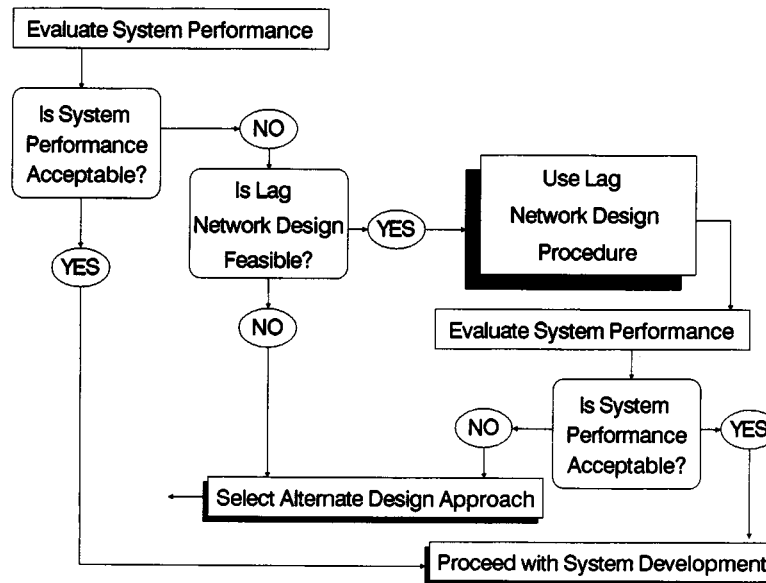**FIGURE 100.14**  A simple electrical lag network.



**FIGURE 100.15**  Life cycle of frequency domain design incorporating lag network compensation.

the crossover frequency and makes the system slower in terms of its transient response. Of course, this would be advantageous if high-frequency noise is present and we wish to reduce its effect. The lag network is an entirely passive device and thus is more economical to instrument than the lead network.

In lead network compensation we actually insert phase lead in the vicinity of the crossover frequency to increase the PM. Thus we realize a specified PM without lowering the medium-frequency system gain. We see that the disadvantages of the lag network are the advantages of the lead network and the advantages of the lag network are the disadvantages of the lead network.

We can attempt to combine the lag network with the lead network into an all-passive structure called a lag-lead network. Generally we obtain better results than we can achieve using either a lead or a lag network. We will consider design using lag-lead networks in our next subsection as well as more complex composite equalization networks.

## Composite Equalizers

In the previous subsection we examined the simplest forms of series equalization: gain adjustment, lead network compensation, and lag network compensation. In this subsection we will consider more complex design examples in which composite equalizers will be used for series compensation. The same design principles used earlier in this section will be used here as well.

### Lag-Lead Network Design

The prime purpose of a lead network is to add phase lead near the crossover frequency to increase the PM. Accompanied with this phase lead is a gain increase that will increase the crossover frequency. This will sometimes cause difficulties if there is much phase lag in the uncompensated system at high frequencies. There may be situations where use of a phase-lead network to achieve a given PM is not possible due to too many high-frequency poles.

The basic idea behind lag network design is to reduce the gain at "middle" frequencies such as to reduce the crossover frequency to a lower value than for the uncompensated system. If the phase lag is less at this lower frequency, then the PM will be increased by use of the lag network. We have seen that is not possible to use a lag network in situations in which there is not a frequency where an acceptable PM would exist if this frequency were the crossover frequency. Even if use of a lag network is possible, the significantly reduced crossover frequency resulting from its use may make the system so slow and sluggish in response to an input that system performance is unacceptable even though the relative stability of the system is acceptable.

Examination of these characteristics or attributes of lead network and lag network compensation suggests that it might be possible to combine the two approaches to achieve the desirable features of each approach. Thus we will attempt to provide attenuation below the crossover frequency to decrease the phase lag at crossover and phase lead closer to the crossover frequency in order to increase the phase lead of the uncompensated system at the crossover frequency.

The transfer function of the basic lag-lead network is

$$G_c(s) = \frac{(1 + s/\omega_2)(1 + s/\omega_3)}{(1 + s/\omega_1)(1 + s/\omega_4)}$$

where $\omega_4 > \omega_3 > \omega_2 > \omega_1$. Often it is desirable that $\omega_2\omega_3 = \omega_1\omega_4$ such that the high-frequency gain of the equalizer is unity. It is generally not desirable that $\omega_1\omega_4 > \omega_2\omega_3$ as this indicates a high-frequency gain greater than 1, and this will require an active network, or gain, and a passive equalizer. It is a fact that we should always be able to realize a linear minimum phase network using passive components only if the network has a rational transfer function with a gain magnitude that is no greater than 1 at any real frequency.

Figure 100.16 illustrates the gain magnitude and phase shift curves for a single-stage lag-lead network equalizer or compensator transfer function. Figure 100.17 illustrates an electrical network realization of a passive lag-lead network equalizer. Parameter matching can be used to determine the electrical network parameters that yield a specified transfer function. Because the relationships between the break frequencies and the equalizer component values are complex, it may be desirable, particularly in preliminary instrumentation of the control system, to use analog or digital computer programming techniques to construct the equalizer. Traditionally, there has been much analog computer simulation of control systems. The more contemporary approach suggests use of digital computer approaches that require numerical approximation of continuous-time physical systems.

Figure 100.18 presents a flowchart that we may use for lag-lead network design. We see that this flowchart has much in common with the charts and design procedures for lead network and lag network design and that each of these approaches first involves determining or obtaining a set of desired specifications for the control system. Next, the form of a trial compensating network and the number of break frequencies in the network are selected. We must then obtain a number of equations, equal to the number of network break frequencies plus 1. One of these equations shows that the gain magnitude is 1 at the crossover frequency. The second equation will be an equation for the phase shift at crossover. It is generally desirable that there be at least two unspecified compensating network break frequencies such that we may use a third equation, the optimality of the phase shift at crossover equation, in which we set $d\beta/d\omega|_{\omega=\omega_c} = 0$. If other equations are needed to represent the design situation, we obtain these from the design specifications themselves.

### General Bode Diagram Design

Figure 100.19 presents a flowchart of a general design procedure for Bode diagram design. As we will see in the next subsection, a minor modification of this flowchart can be used to accomplish design using minor-loop feedback
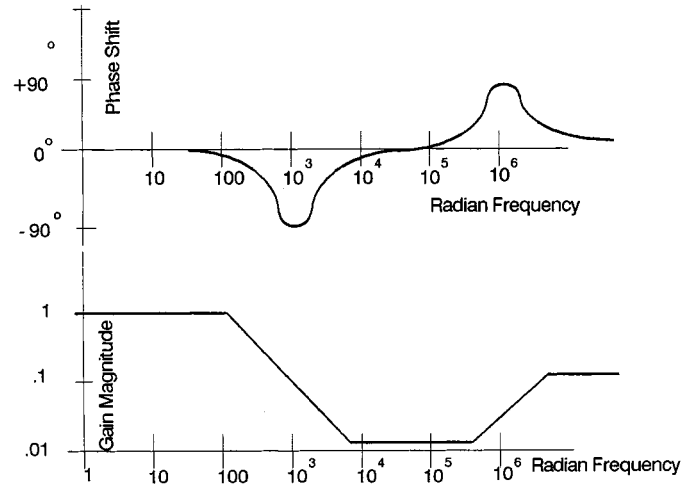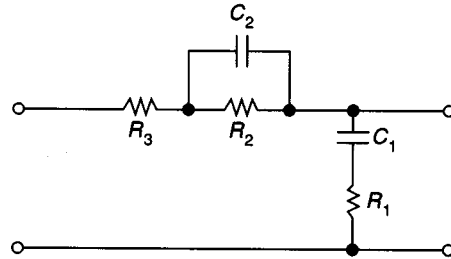
**FIGURE 100.16** Phase shift and gain curves for a simple lag-lead network.



$$G_c(s) = \frac{(1 + s/\omega_2)(1 + s/\omega_3)}{(1 + s/\omega_1)(1 + s/\omega_4)}$$

$$= \frac{(1 + R_1C_1s)(1 + R_2C_2s)}{1 + (R_1C_1 + R_2C_1 + R_3C_1 + R_2C_2)s + (R_1R_2C_1C_2 + R_2R_3C_1C_2)s^2}$$

Special case: $R_3 = 0$

$$\omega_2 = \frac{1}{R_1C_1}, \quad \omega_3 = \frac{1}{R_2C_2}, \quad \omega_1\omega_4 = \omega_2\omega_3, \quad \omega_1 + \omega_4 = \omega_2 + \omega_3 + \frac{1}{R_1C_2}$$

**FIGURE 100.17** Simple electrical lag-lead network.

or a combination of minor-loop and series equations. These detailed flowcharts for Bode diagram design are, of course, part of the overall design procedure of Fig. 100.9.

Much experience leads to the conclusion that satisfactory linear systems control design using frequency response approaches is such that the crossover frequency occurs on a gain magnitude curve which has a −1 slope at the crossover frequency. In the vicinity of crossover we may approximate any minimum phase transfer function, with crossover on a −1 slope, by

$$G(s) = G_f(s)G_c(s) = \frac{\omega_c\omega_1^{n-1}(1 + s/\omega_1)^{n-1}}{s^n(1 + s/\omega_2)^{m-1}} \quad \text{for } \omega_1 > \omega_c > \omega_2$$
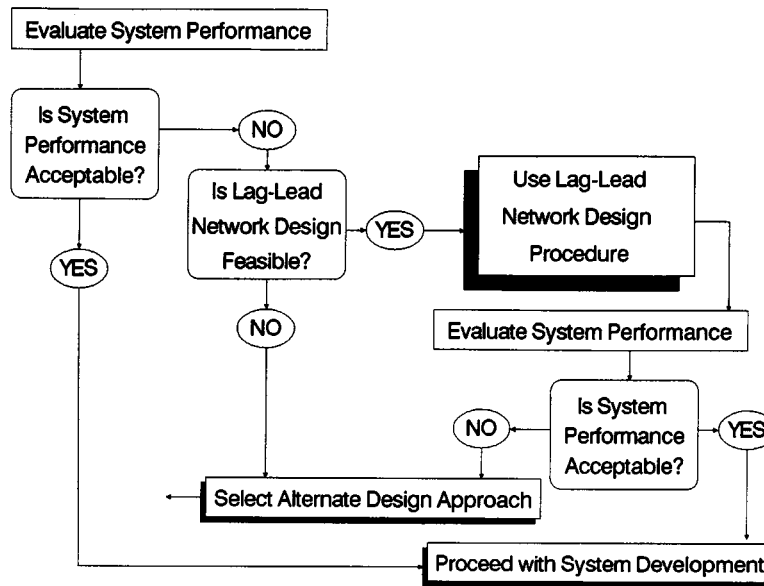
**FIGURE 100.18**  Life cycle of frequency domain design incorporating lag-lead network compensation.



**FIGURE 100.19**  Life cycle of frequency domain design incorporating general Bode diagram compensation procedure.

Here $\omega_1$ is the break frequency just prior to crossover and $\omega_2$ is the break frequency just after crossover. It is easy to verify that we have $|G(j\omega_c)| = 1$ if $\omega_1 > \omega_c > \omega_2$. Figure 100.20 illustrates this rather general approximation to a compensated system Bode diagram in the vicinity of the crossover frequency. We will conclude this subsection by determining some general design requirements for a system with this transfer function and the associated Bode asymptotic gain magnitude diagram of Fig. 100.20.

There are three unknown frequencies in the foregoing equation. Thus we need three requirements or equations to determine design parameters. We will use the same three requirements used thus far in all our efforts in this section, namely:

---

**FIGURE 100.20** Illustration of generic gain magnitude in the vicinity of crossover.

1. The gain at crossover is 1.
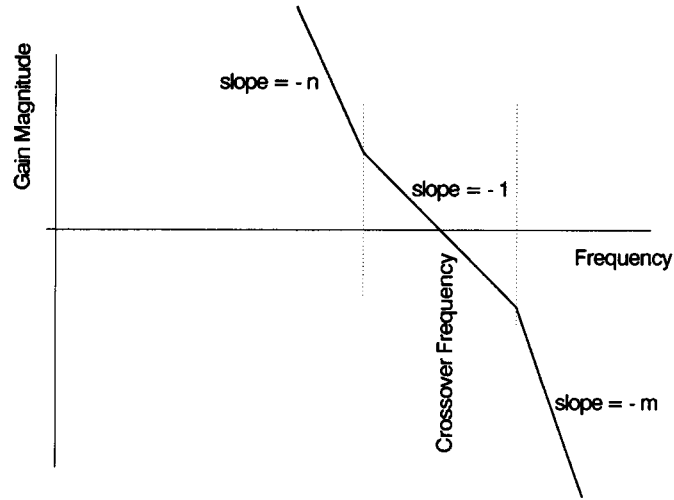2. The PM is some specified value.
3. The PM at crossover is the maximum possible for a given $\omega_2/\omega_1$ ratio.

We see that the first requirement, that the gain is 1 at the crossover frequency, is satisfied by the foregoing equation if the crossover frequency occurs on the −1 slope portion of the gain curve as assumed in Fig. 100.20. We use the arctangent approximation to obtain the phase shift in the vicinity of crossover as

$$\beta(\omega) = -\frac{n\pi}{2} + (n-1)\left(\frac{\pi}{2} - \frac{\omega_1}{\omega}\right) - (m-1)\frac{\omega}{\omega_2}$$

To satisfy requirement 3 we set

$$\frac{d\beta(\omega)}{d\omega}\bigg|_{\omega=\omega_c} = 0 = \frac{(n-1)\omega_1}{\omega_c^2} - \frac{m-1}{\omega_2}$$

and obtain

$$\omega_c^2 = \frac{n-1}{m-1}\,\omega_1\,\omega_2$$

as the optimum setting for the crossover frequency. Substitution of the "optimum" frequency given by the foregoing into the phase shift equation results in

$$\beta(\omega_c) = \frac{-\pi}{2} - 2\sqrt{(m-1)(n-1)}\sqrt{\frac{\omega_1}{\omega_2}}$$

We desire a specific PM here, and so the equalizer break frequency locations are specified. There is a single parameter here that is unspecified, and an additional equation must be found in any specific application. Alternately, we could simply assume a nominal crossover frequency of unity or simply normalize frequencies $\omega_1$ and $\omega_2$ in terms of the crossover frequency by use of the normalized frequencies $\omega_1 = W_1\omega_c$ and $\omega_2 = W_2\omega_c$.
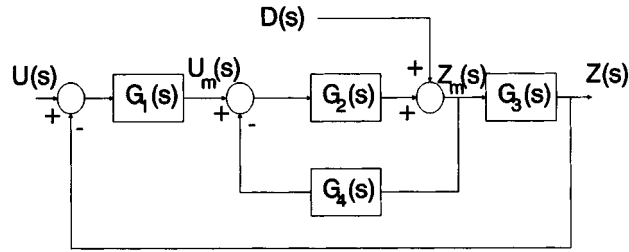
**FIGURE 100.21** Feedback control system with a single minor loop and output disturbance.

It is a relatively simple matter to show that for a specified PM expressed in radians, we obtain for the normalized break frequencies

$$W_1 = \frac{\omega_1}{\omega_c} = \frac{PM}{2(n-1)}$$

$$W_2 = \frac{\omega_2}{\omega_c} = \frac{2(m-1)}{PM}$$

It is relatively easy to implement this suggested Bode diagram design procedure which is based upon considering only these break frequencies immediately above and below crossover and which approximate all others. Break frequencies far below crossover are approximated by integrations or differentiations, that is, poles or zeros at $s = 0$, and break frequencies far above the crossover frequency are ignored.

## Minor-Loop Design

In our efforts thus far in this section we have assumed that compensating networks would be placed in series with the fixed plant and then a unity feedback ratio loop closed around these elements to yield the closed-loop system. In many applications it may be physically convenient, perhaps due to instrumentation considerations, to use one or more minor loops to obtain a desired compensation of a fixed plant transfer function.

For a single-input–single-output linear system there are no theoretical advantages whatever to any minor-loop compensation to series compensation as the same closed-loop transfer function can be realized by all procedures. However, when there are multiple inputs or outputs, then there may be considerable advantages to minor-loop design as contrasted to series compensation design. Multiple inputs often occur when there is a single-signal input and one or more noise or disturbance inputs present and a task of the system is to pass the signal inputs and reject the noise inputs. Also there may be saturation-type nonlinearities present, and we may be concerned not only with the primary system output but also with keeping the output at the saturation point within bounds such that the system remains linear. Thus there are reasons why minor-loop design may be preferable to series equalization.

We have discussed block diagrams elsewhere in this handbook. It is desirable here to review some concepts that will be of value for our discussion of minor-loop design. Figure 100.21 illustrates a relatively general linear control system with a single minor loop. This block diagram could represent many simple control systems. $G_1(s)$ could represent a discriminator and series compensation and $G_2(s)$ could represent an amplifier and that part of a motor transfer function excluding the final integration to convert velocity to position. $G_3(s)$ might then represent an integrator. $G_4(s)$ would then represent a minor-loop compensation transfer function, such as that of a tachometer.

The closed-loop transfer function for this system is given by

$$\frac{Z(s)}{U(s)} = H(s) = \frac{G_1(s)G_2(s)G_3(s)}{1 + G_2(s)G_4(s) + G_1(s)G_2(s)G_3(s)}$$

It is convenient to define several other transfer functions that are based on the block diagram in Fig. 100.21. First there is the minor-loop gain

$$G_m(s) = G_2(s)G_4(s)$$

which is just the loop gain of the minor loop only. The minor loop has the transfer function

$$\frac{Z_m(s)}{U_m(s)} = H_m(s) = \frac{G_2(s)}{1 + G_2(s)G_4(s)} = \frac{G_2(s)}{1 + G_m(s)}$$

There will usually be a range or ranges of frequency for which the minor-loop gain magnitude is much less than 1, and we then have

$$\frac{Z_m(s)}{U_m(s)} = H_m(s) \approx G_2(s) \qquad |G_m(\omega)| << 1$$

There will also generally be ranges of frequency for which the minor-loop gain magnitude is much greater than 1, and we then have

$$\frac{Z_m(s)}{U_m(s)} = H_m(s) \approx \frac{1}{G_4(s)} \qquad |G_m(\omega)| >> 1$$

We may use these two relations to considerably simplify our approach to the minor-loop design problem. For illustrative purposes, we will use two major-loop gain functions. First we will consider the major-loop gain with the minor-loop-compensating network removed such that $G_4(s) = 0$. This represents the standard situation we have examined in the last subsection. This uncompensated major-loop transfer function is

$$G_{Mu}(s) = G_1(s)G_2(s)G_3(s)$$

With the minor-loop compensation inserted, the major-loop gain, the input-output transfer function with the unity ratio feedback open, is

$$G_{Mc}(s) = \frac{G_1(s)G_2(s)G_3(s)}{1 + G_m(s)}$$

We may express the complete closed-loop transfer function in the form

$$\frac{Z(s)}{U(s)} = H(s) = \frac{G_{Mc}(s)}{1 + G_{Mc}(s)}$$

A particularly useful relationship may be obtained by combining the last three equations into one equation of the form

$$G_{Mc}(s) = \frac{G_{Mu}(s)}{1 + G_m(s)}$$

We may give this latter equation a particularly simple interpretation. For frequencies where the minor-loop gain $G_m(s)$ is low, the minor-loop–closed major-loop transfer function $G_{Mc}(s)$ is approximately that of the minor-loop–open major-loop transfer function $G_{Mu}$ in that

$$G_{Mc}(s) \approx G_{Mu}(s) \qquad \left| G_m(\omega) \right| << 1$$

For frequencies where the minor-loop gain $G_m(s)$ is high, the minor-loop–closed major-loop transfer function is just

$$G_{Mc}(s) \approx \frac{G_{Mu}(s)}{G_m(s)} \qquad \left| G_m(\omega) \right| >> 1$$

This has an especially simple interpretation on the logarithmic frequency plots we use for Bode diagrams for we may simply subtract the minor-loop gain $G_m(s)$ from the minor-loop–open major-loop gain $G_{Mu}(s)$ to obtain the compensated system gain as the transfer function $G_{Mc}(s)$.

The last several equations are the key relations for minor-loop design using this frequency response approach. These relations indicate that some forms of series compensation yield a given major-loop transfer function $G_{Mc}(s)$ which will not be appropriate for realization by minor-loop compensation. In particular, a lead network series compensation cannot be realized by means of equivalent minor-loop compensation. The gain of the fixed plant $G_{Mu}(s)$ is raised at high frequencies due to the use of a lead network compensation. Also, we see that $G_{Mc}(s)$ can only be lowered by use of a minor-loop gain $G_m(s)$.

A lag network used for series compensation will result in a reduction in the fixed plant gain $\left| G_{Mu}(\omega) \right|$ at all high frequencies. This can only be achieved if the minor-loop transfer gain $G_m(s)$ is constant for high frequencies. In some cases this may be achievable but often will not be. It is possible to realize the equivalent of lag network series equalization by means of a minor-loop equalization for systems where the low- and high-frequency behavior of $G_{Mu}(s)$, or $G_f(s)$, and $G_{Mc}(s)$ are the same and where the gain magnitude of the compensated system $\left| G_{Mc}(s) \right|$ is at no frequency any greater than is the gain magnitude of the fixed plant $\left| G_f(s) \right|$ or the minor-loop–open major-loop transfer function $\left| G_{Mu}(s) \right|$. Thus we see that lag-lead network series equalization is an ideal type of equalization to realize by means of equivalent minor-loop equalization. Figures 100.9 and 100.19 represent flowcharts of a suggested general design procedure for minor-loop compensator design as well as for the series equalization approaches we examined previously.

In our work thus far we have assumed that parameters were constant and known. Such is seldom the case, and we must naturally be concerned with the effects of parameter variations, disturbances, and nonlinearities upon system performance. Suppose, for example, that we design a system with a certain gain assumed as $K_1$. If the system operates open loop and the gain $K_1$ is in cascade or series with the other input-output components, then the overall transfer function changes by precisely the same factor as $K_1$ changes. If we have an amplifier with unity ratio feedback around a gain $K_1$, the situation is much different. The closed-loop gain would nominally be $K_1/(1 + K_1)$, and a change to $2K_1$ would give a closed-loop gain $2K_1/(1 + 2K_1)$. If $K_1$ is large, say $10_3$, then the new gain is 0.99950025, which is a percentage change of less than 0.05% for a change in gain of 100%.

Another advantage of minor-loop feedback occurs when there are output disturbances such as those due to wind gusts on an antenna. We consider the system illustrated in Fig. 100.21. The response due to $D(s)$ alone is

$$\frac{Z(s)}{D(s)} = \frac{1}{1 + G_2(s)\, G_4(s) + G_1(s)\, G_2(s)}$$

When we use the relation for the minor-loop gain

$$G_m(s) = G_2(s)\, G_4(s)$$

and the major-loop gain

$$G_{Mc}(s) = \frac{G_1(s)G_2(s)}{1 + G_2(s)G_4(s)}$$

we can rewrite the response due to $D(s)$ as

$$\frac{Z(s)}{D(s)} = \frac{1}{[1 + G_m(s)]G_{Mc}(s)}$$

Over the range of frequency where $|G_{Mc}(j\omega)| \gg 1$, such that the corrected loop gain is large, the attenuation of a load disturbance is proportional to the uncorrected loop gain. This is generally larger over a wider frequency range than the corrected loop gain magnitude $|G_{Mc}(j\omega)|$, which is what the attenuation would be if series compensation were used.

Over the range of frequencies where the minor-loop gain is large but where the corrected loop gain is small, that is, where $|G_m(j\omega)| > 1$ and $|G_{Mc}(j\omega)| < 1$, we obtain for the approximate response due to the disturbance

$$\frac{Z(s)}{D(s)} \approx G_m(s)$$

and the output disturbance is therefore seen to be attenuated by the minor-loop gain rather than unattenuated as would be the case if series compensation had been used. This is, of course, highly desirable.

At frequencies where both the minor-loop gain transfer and the major-loop gain are small we have $Z(s)/D(s) \approx 1$, and over this range of frequencies neither minor-loop compensation nor series equalization is useful in reducing the effect of a load disturbance. Thus, we have shown here that there are quite a number of advantages to minor-loop compensation as compared to series equalization. Of course, there are limitations as well.

## Summary

In this section, we have examined the subject of linear system compensation by means of the frequency response method of Bode diagrams. Our approach has been entirely in the frequency domain. We have discussed a variety of compensation networks, including:

1. Gain attenuation
2. Lead networks
3. Lag networks
4. Lag-lead networks and composite equalizers
5. Minor-loop feedback

Despite its age, the frequency domain design approach represents a most useful approach for the design of linear control systems. It has been tested and proven in a great many practical design situations.

## Defining Terms

**Bode diagram:** A graph of the gain magnitude and frequency response of a linear circuit or system, generally plotted on log-log coordinates. A major advantage of Bode diagrams is that the gain magnitude plot will look like straight lines or be asymptotic to straight lines. H.W. Bode, a well-known Bell Telephone Laboratories researcher, published *Network Analysis and Feedback Amplifier Design* in 1945. The approach, first described there, has been refined by a number of other workers over the past half-century.

**Crossover frequency:** The frequency where the magnitude of the open-loop gain is 1.

**Equalizer:** A network inserted into a system that has a transfer function or frequency response designed to compensate for undesired amplitude, phase, and frequency characteristics of the initial system. Filter and equalizer are generally synonymous terms.

**Lag network:** In a simple phase-lag network, the phase angle associated with the input-output transfer function is always negative, or lagging. Figures 100.13 and 100.14 illustrate the essential characteristics of a lag network.

**Lag-lead network:** The phase shift versus frequency curve in a phase lag-lead network is negative, or lagging, for low frequencies and positive, or leading, for high frequencies. The phase angle associated with the input-output transfer function is always positive, or leading. Figures 100.16 and 100.17 illustrate the essential characteristics of a lag-lead network, or composite equalizer.

**Lead network:** In a simple phase-lead network, the phase angle associated with the input-output transfer function is always positive, or leading. Figures 100.10 and 100.11 illustrate the essential characteristics of a lead network.

**Series equalizer:** In a single-loop feedback system, a series equalizer is placed in the single loop, generally at a point along the forward path from input to output where the equalizer itself consumes only a small amount of energy. In Fig. 100.21, $G_1(s)$ could represent a series equalizer. $G_1(s)$ could also be a series equalizer if $G_4(s) = 0$.

**Specification:** A statement of the design or development requirements to be satisfied by a system or product.

**Systems engineering:** An approach to the overall life cycle evolution of a product or system. Generally, the systems engineering process comprises a number of phases. There are three essential phases in any systems engineering life cycle: formulation of requirements and specifications, design and development of the system or product, and deployment of the system. Each of these three basic phases may be further expanded into a larger number. For example, deployment generally comprises operational test and evaluation, maintenance over an extended operational life of the system, and modification and retrofit (or replacement) to meet new and evolving user needs.

## Related Topic

11.1 Introduction

## References

J.L. Bower and P.M. Schultheiss, *Introduction to the Design of Servomechanisms*, New York: Wiley, 1958.

A.P. Sage, *Linear Systems Control,* Champaign, Ill.: Matrix Press, 1978.

A.P. Sage, *Systems Engineering,* New York: Wiley, 1992.

M.G. Singh, Ed., *Systems and Control Encyclopedia,* Oxford: Pergamon, 1987.

## Further Information

Many of the practical design situations used to test the frequency domain design approach are described in the excellent classic text by Bower and Schultheiss [1958]. A rather detailed discussion of the approach may also be found in Sage [1978] on which this discussion is, in part, based. A great variety of control systems design approaches, including frequency domain design approaches, are discussed in a recent definitive control systems encyclopedia [Singh, 1987], and there are a plethora of new introductory control systems textbooks that discuss it as well. As noted earlier, frequency domain design, in particular, and control systems design, in general, constitute one facet of systems engineering effort, such as described in Sage [1992].

# 100.4   Root Locus

*Benjamin C. Kuo*

**Root locus** represents a trajectory of the roots of an algebraic equation with constant coefficients when a parameter varies. The technique is used extensively for the analysis and design of linear time-invariant control

systems. For linear time-invariant control systems the roots of the characteristic equation determine the stability of the system. For a stable continuous-data system the roots must all lie in the left half of the $s$ plane. For a digital control system to be stable, the roots of the characteristic equation must all lie inside the unit circle $|z| = 1$ in the $z$ plane. Thus, in the $s$ plane, the imaginary axis is the stability boundary, whereas in the z plane the stability boundary is the unit circle. The location of the characteristic equation roots with respect to the stability boundary also determine the relative stability, i.e., the degree of stability, of the system.

For a linear time-invariant system with continuous data, the characteristic equation can be written as

$$F(s) = P(s) + KQ(s) = 0 \tag{100.33}$$

where $P(s)$ is an $N$th-order polynomial of $s$,

$$P(s) = s^N + a_1 s^{N-1} + \ldots + a_{N-1}s + a_N \tag{100.34}$$

and $Q(s)$ is the $M$th-order polynomial of $s$,

$$Q(s) = s^M + b_1 s^{M-1} + \ldots + b_{M-1}s + b_M \tag{100.35}$$

where $N$ and $M$ are positive integers. The real constant $K$ can vary from $-\infty$ to $+\infty$. The coefficients $a_1$, $a_2$, …, $a_N$, $b_1$, $b_2$, …, $b_M$ are real. As $K$ is varied from $-\infty$ to $+\infty$, the roots of Eq. (100.33) trace out continuous trajectories in the $s$ plane called the *root loci*.

The above development can be extended to digital control systems by replacing $s$ with $z$ in Eqs. (100.33) through (100.35).

## Root Locus Properties

The root locus problem can be formulated from Eq. (100.33) by dividing both sides of the equation by the terms that do not contain the variable parameter $K$. The result is

$$1 + \frac{KQ(s)}{P(s)} = 0 \tag{100.36}$$

For a closed-loop control system with the loop transfer function $KG(s)H(s)$, where the gain factor $K$ has been factored out, the characteristic equation is known to be the zeros of the rational function

$$1 + KG(s)H(s) = 0 \tag{100.37}$$

Since Eqs. (100.36) and (100.37) have the same form, the general root locus problem can be formulated using Eq. (100.36).

Equation (100.37) is written

$$G(s)H(s) = -\frac{1}{K} \tag{100.38}$$

To satisfy the last equation, the following conditions must be met simultaneously:

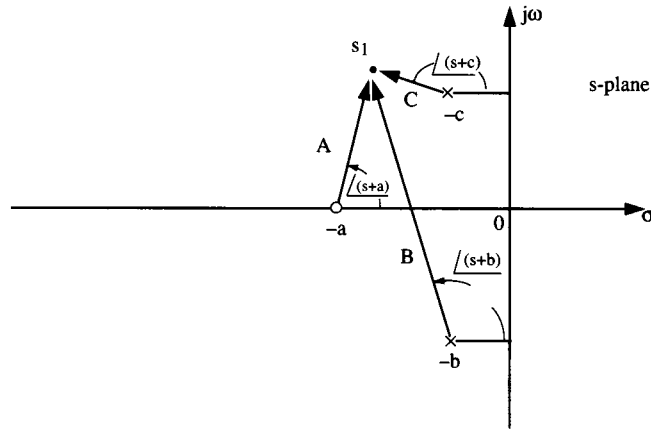$$\text{Condition on magnitude:} \quad |G(s)H(s)| = \frac{1}{|K|} \tag{100.39}$$

**FIGURE 100.22**  Graphical interpretation of magnitude and angle conditions of root loci.

where $K$ varies between $-\infty$ and $+\infty$.

$$\text{Conditions on angles: } \angle G(s)H(s) = (2k + 1)\pi \quad K \geq 0$$

$$= \text{odd multiples of } \pi \text{ rad} \qquad (100.40)$$

$$\angle G(s)H(s) = 2k\pi \quad K \leq 0$$

$$= \text{even multiples of } \pi \text{ rad} \qquad (100.41)$$

where $k = 0, \pm 1, \pm 2, \ldots, \pm$ any integer.

In general, the conditions on angles in Eqs. (100.40) and (100.41) are used for the construction of the root loci, whereas the condition on magnitude in Eq. (100.39) is used to find the value of $K$ on the loci once the loci are drawn. Let $KG(s)H(s)$ be of the form

$$KG(s)H(s) = \frac{K(s + a)}{(s + b)(s + c)} \qquad (100.42)$$

Applying Eqs. (100.40) and (100.41) to the last equation, the angles conditions are

$$K \geq 0: \quad \angle G(s)H(s) = \angle(s + a) - \angle(s + b) - \angle(s + c)$$

$$= (2k + 1)\pi \qquad (100.43)$$

$$K \leq 0: \quad \angle G(s)H(s) = \angle(s + a) - \angle(s + b) - \angle(s + c)$$

$$= 2k\pi \qquad (100.44)$$

where $k = 0, \pm 1, \pm 2, \ldots$ . The graphical interpretation of the last two equations is shown in Fig. 100.22. For the point $s_1$ to be a point on the root locus, the angles of the phasors drawn from the poles and zeros of $G(s)H(s)$ to $s_1$ must satisfy Eq. (100.43) or (100.44) depending on the sign of $K$. Applying the magnitude condition of Eq. (100.39) to (100.42), the magnitude of $K$ is expressed as

$$|K| = \frac{|s + b||s + c|}{|s + a|} = \frac{B \cdot C}{A} \qquad (100.45)$$

where *A, B,* and *C* are the lengths of the phasors drawn from the poles and zeros of $G(s)H(s)$ to the point $s_1$.

The following properties of the root loci are useful for sketching the root loci based on the pole-zero configuration of $G(s)H(s)$. Many computer programs, such as the ROOTLOCI in the ACSP software package [Kuo, 1991b], are available for computing and plotting the root loci. The proofs and derivations of these properties can be carried out from Eqs. (100.39), (100.40), and (100.41) [Kuo, 1991a].

***Starting Points (K 5 0 Points).*** The points at which $K = 0$ on the root loci are at the poles of $G(s)H(s)$.

***Ending Points (K 56` Points).*** The points at which $K = \pm\infty$ on the root loci are at the zeros of $G(s)H(s)$. The poles and zeros referred to above include those at s = ∞.

***Number of Root Loci.*** The total number of root loci of Eq. (100.37) equals the higher of the number of poles and zeros of $G(s)H(s)$.

***Symmetry of Root Loci.*** The root loci are symmetrical with respect to the axes of symmetry of the pole-zero configuration of $G(s)H(s)$. In general, the root loci are symmetrical at least to the real axis of the complex $s$ plane.

***Asymptotes of the Root Loci.*** Asymptotes of the root loci refer to the behavior of the root loci at $|s| = \infty$ when the number of poles and zeros of $G(s)H(s)$ is not equal. Let N denote the number of finite poles of $G(s)H(s)$ and M be the number of finite zeros of $G(s)H(s)$. In general, $2|N - M|$ of the loci will approach infinity in the s plane. The properties of the root loci at $|s| = \infty$ are described by the angles and the intersects of the asymptotes. When $N \neq M$, the angles of the asymptotes are given by

$$\phi_k = \begin{cases} \dfrac{(2k + 1)\pi}{|N - M|} & K \geq 0 \qquad (100.46) \\[4mm] \dfrac{2k\,\pi}{|N - M|} & K \leq 0 \qquad (100.47) \end{cases}$$

where $k = 0, 1, 2, \ldots, |N - M| - 1$.

The asymptotes intersect on the real axis at

$$\sigma = \frac{\sum \text{finite poles of } G(s)H(s) - \sum \text{finite zeros of } G(s)H(s)}{N - M} \qquad (100.48)$$

***Root Loci on the Real Axis.*** The entire real axis of the $s$ plane is occupied by the root loci. When $K > 0$, root loci are found in sections of the real axis to the right of which the total number of poles and zeros of $G(s)H(s)$ is *odd*. When $K < 0$, root loci are found in sections to the right of which the total number of poles and zeros of $G(s)H(s)$ is *even*.

As a summary of the root locus properties discussed above, the properties of the root loci of the following equation are displayed in Fig. 100.23.

$$s^3 + 2s^2 + 2s + K(s + 3) = 0 \qquad (100.49)$$

Dividing both sides of the last equation by the terms that do not contain *K* we get

$$1 + KG(s)H(s) = 1 + \frac{K(s + 3)}{s(s^2 + 2s + 2)} \qquad (100.50)$$

Thus, the poles of $G(s)H(s)$ are at $s = 0$, $s = -1 + j$, and $s = -1 - j$. The zero of $G(s)H(s)$ is at $z = -3$.

As shown in Fig. 100.23, the $K = 0$ points on the root loci are at the poles of $G(s)H(s)$, and the $K = \pm\infty$ points are at the zeros of $G(s)H(s)$. Since $G(s)H(s)$ has two zeros at $s = \infty$, two of the three root loci approach
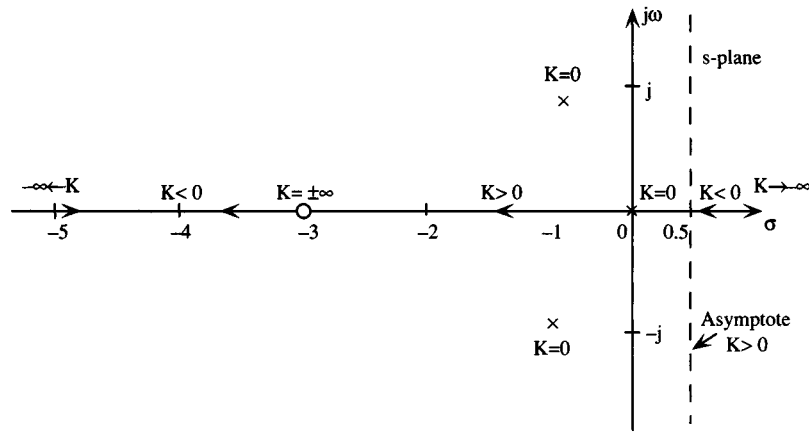
**FIGURE 100.23**  Some properties of the root loci of $G(s)H(s) = K(s + 3)/[s(s^2 + 2s + 2)]$.

infinity in the $s$ plane. The root loci are symmetrical to the real axis of the $s$ plane, since the pole-zero configuration of $G(s)H(s)$ is symmetrical to the axis. The asymptotes of the two root loci that approach infinity are characterized by Eqs. (100.46) through (100.48). The angles of the asymptotes are:

$$K \geq 0: \quad \phi_k = \frac{(2k + 1)\pi}{3 - 1} \qquad k = 0, 1 \tag{100.51}$$

$$K \leq 0: \quad \phi_k = \frac{2k\pi}{3 - 1} \qquad k = 0, 1 \tag{100.52}$$

Thus, for $K \geq 0$, the angles of the asymptotes are $\phi_0 = 90°$ and $\phi_1 = 270°$. For $K \leq 0$, $\phi_0 = 0°$ and $\phi_1 = 180°$.
  The intersect of the asymptotes is at

$$\sigma = \frac{-1 + j - 1 - j - (-3)}{3 - 1} = \frac{1}{2} \tag{100.53}$$

The root loci on the real axis are as indicated in Fig. 100.23.

***Angles of Departure and Arrival.***    The slope of the root locus in the vicinity of a pole of $G(s)H(s)$ is measured at the *angle of departure* and that in the vicinity of a zero of $G(s)H(s)$ is measured at the *angle of arrival*.
  The angle of departure (arrival) of a root locus at a pole (zero) of $G(s)H(s)$ is determined by assigning a point $s_1$ to the root locus that is very close to the pole (zero) and applying the angle conditions of Eqs. (100.40) or (100.41). Figure 100.24 illustrates the calculation of the angles of arrival and departure of the root locus at the pole $s = -1 + j$. We assign a point $s_1$ that is on the locus for $K > 0$ near the pole and draw phasors from *all* the poles and the zero of $G(s)H(s)$ to this point. The angles made by the phasors with respect to the real axis must satisfy the angle condition in Eq. (100.46). Let the angle of the phasor drawn from $-1 + j$ to $s_1$ be designated as $\theta$, which is the angle of departure; the angles drawn from the other poles and zero can be approximated by regarding $s_1$ as being very close to $-1 + j$. Thus, Eq. (100.46) leads to

$$\angle G(s_1)H(s_1) = -\theta - 135° - 90° + 26.6° = -180° \tag{100.54}$$

or $\theta = -18.4°$. For the angle of arrival of the root locus at the pole $s = -1 + j$, we assign a point $s_1$ on the root loci for $K < 0$ near the pole. Drawing phasors from all the poles and the zero of $G(s)H(s)$ to $s_1$ and applying the angle condition in Eq. (100.47), we have
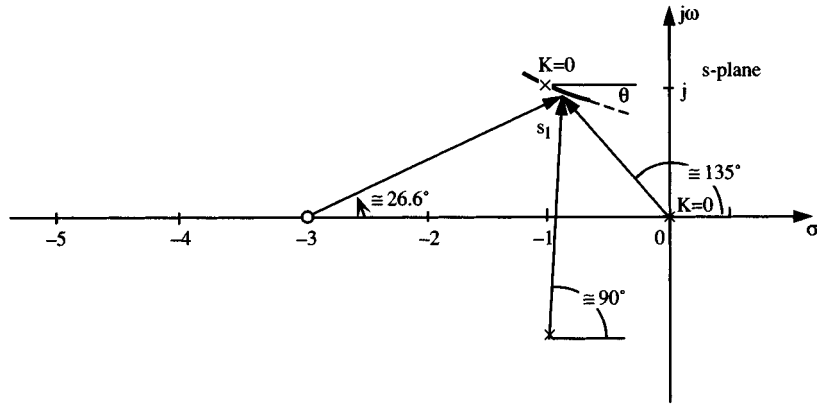
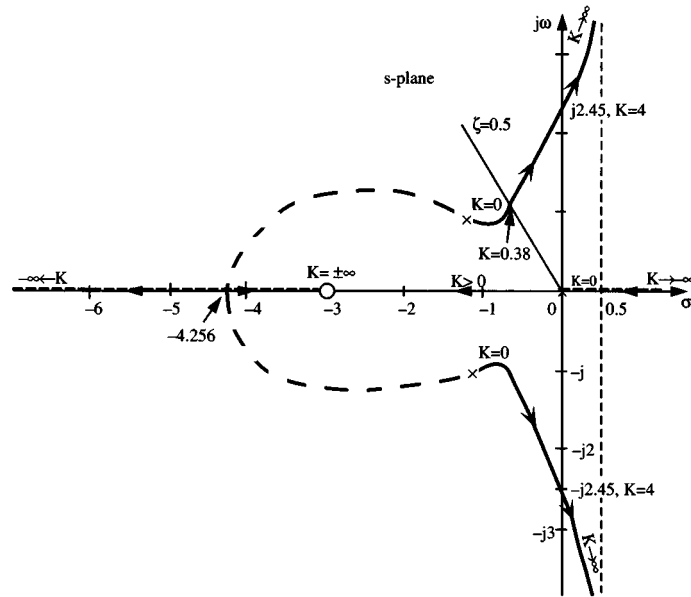**FIGURE 100.24** Angle of arrival and departure calculations.



**FIGURE 100.25** The complete root loci of $G(s)H(s) = K(s + 3)/[s(s^2 + 2s + 2)]$.

$$\angle G(s_1)H(s_1) = -\theta - 135° - 90° + 26.6° = 0° \qquad (100.55)$$

Thus, the angle of arrival of the root locus for $K < 0$ is $\theta = 198.4°$. Similarly, we can show that the angles of arrival and departure of the root locus at $s = -3$ are $180°$ and $0°$, respectively.

***Intersection of the Root Loci with the Imaginary Axis.***    The points where the root loci intersect the imaginary axis (if there is any) in the $s$ plane, and the corresponding values of $K$, may be determined by means of the Routh-Hurwitz stability criterion [Kuo, 1991a]. The root locus program can also be used on a computer to give the intersects.

The complete root loci in Fig. 100.25 show that the root loci intersect the imaginary axis at $s = \pm j2.45$, and the value of $K$ is 4. The system is stable for $0 \le K < 4$.

***Breakaway Points of the Root Loci.***    Breakaway points on the root loci correspond to multiple-order roots of the equation. At a breakaway point several root loci converge and then break away in different directions. The breakaway point can be real or complex. The latter case must be in complex conjugate pairs.

The breakaway points of the root loci of Eq. (100.37) must satisfy the following condition:

$$\frac{dG(s)H(s)}{ds} = 0 \qquad (100.56)$$

On the other hand, not all solutions of Eq. (100.56) are breakaway points. To satisfy as a breakaway point, the point must also lie on the root loci, or satisfy Eq. (100.37). Applying Eq. (100.56) to the function $G(s)H(s)$ given in Eq. (100.50), we have the equation that the breakaway points must satisfy,

$$2s^3 + 11s^2 + 12s + 6 = 0 \qquad (100.57)$$

The roots of the last equation are $s = -4.256$, $-0.622 + j0.564$ and $-0.622 - j0.564$. As shown in Fig. 100.25, only the solution $s = -4.256$ is a breakaway point on the root loci.

## Root Loci of Digital Control Systems

The root locus analysis presented in the preceding subsections can be applied to digital control systems without modifying the basic principles. For a linear time-invariant digital control system, the transfer functions are expressed in terms of the $z$-transform rather than the Laplace transform. The relationship between the $z$-transform variable $z$ and the Laplace transform variable $s$ is

$$z = e^{Ts} \qquad (100.58)$$

where $T$ is the sampling period in seconds. Typically, the characteristic equation roots are solutions of the equation

$$1 + KGH(z) = 0 \qquad (100.59)$$

where $K$ is the variable gain parameter. The root loci for a digital control system are constructed in the complex $z$ plane. All the properties of the root loci in the $s$ plane apply readily to the $z$ plane, with the exception that the stability boundary is now the unit circle $|z| = 1$. That is, the system is stable if all the characteristic equation roots lie inside the unit circle.

As an illustration, the open-loop transfer function of a digital control system is given as

$$G(z) = \frac{K(z + 0.1)}{z(z - 1)} \qquad (100.60)$$

The characteristic equation of the closed-loop system is

$$z(z - 1) + K(z + 0.1) = 0$$

The root loci of the system are shown in Fig. 100.26. Notice that the system is stable for $0 \le K < 2.22$. When $K = 2.22$, one root is at $z = -1$, which is on the stability boundary.

## Design with Root Locus

The root locus diagram of the characteristic equation of a closed-loop control system can be used for design purposes. The roots of the characteristic equation can be positioned in the $s$ plane (or the $z$ plane for digital control systems) to realize a certain desired relative stability or damping of the system. It should be kept in mind that the zeros of the closed-loop transfer function also affect the relative stability of the system, although the absolute stability is strictly governed by the characteristic equation roots.
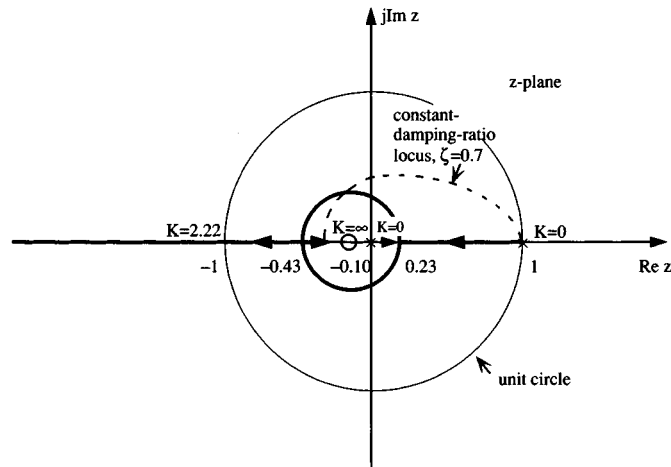
**FIGURE 100.26**  Root loci in the *z* plane for a digital control system.

As an illustrative example, the constant-damping ratio line for $\zeta = 0.5$ is shown in Fig. 100.25. The intersect of the $\zeta = 0.5$ line and the root locus corresponds to $K = 0.38$. Let us assume that we want to keep the relative damping at approximately 0.5 but the gain $K$ should be increased tenfold. The following cascade controller is applied to the system [Evans, 1948]:

$$G_c(s) = \frac{1 + 5s}{1 + 50s} \tag{100.61}$$

The open-loop transfer function of the compensated system is now

$$G_c(s)G(s)H(s) = \frac{0.1K(s + 3)(s + 0.2)}{s(s + 0.02)(s^2 + 2s + 2)} \tag{100.62}$$

Figure 100.27 shows the root locus diagram of the compensated system for $K \geq 0$. The shape of the complex root loci is not appreciably affected by the controller, but the value of $K$ that corresponds to a relative damping ratio of 0.5 is now approximately 3.9.

In a similar manner the root loci of digital control systems can be reshaped in the *z* plane for design. The constant-damping ratio locus in the *z* plane is shown in Fig. 100.26.

## Defining Terms

**Angles of departure and arrival:**  The slope of the root locus in the vicinity of a pole of $G(s)H(s)$ is measured as the angle of departure, and that in the vicinity of a zero of $G(s)H(s)$ is measured as the angle of arrival.

**Asymptotes of root loci:**  The behavior of the root loci at $|s| = \infty$ *when the number of poles and zeros of $G(s)H(s)$ is not equal.*

**Breakaway points of the root loci:**  Breakaway points on the root loci correspond to multiple-order roots of the equation.

**Root locus:**  The trajectory of the roots of an algebraic equation with constant coefficient when a parameter varies.

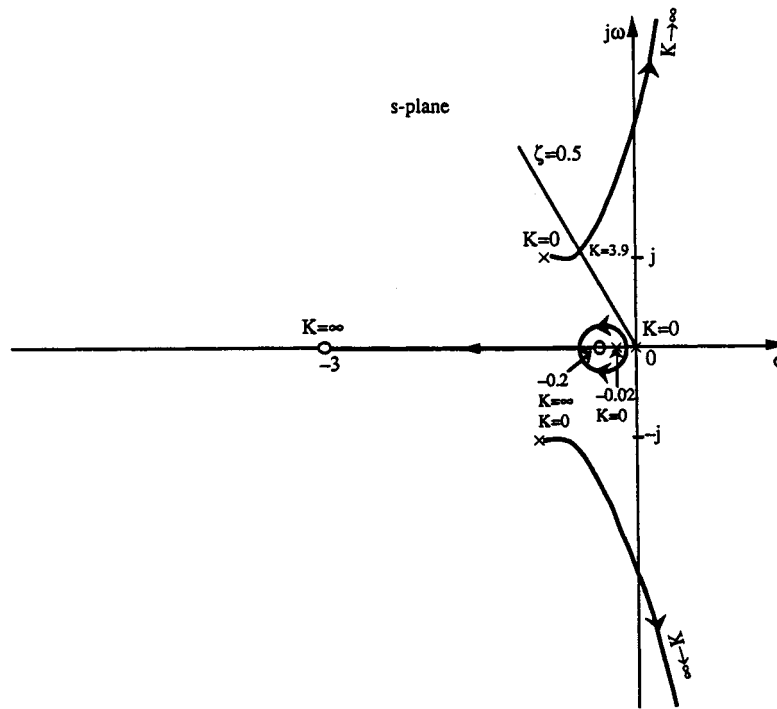## Related Topics

6.1 Definitions and Properties • 12.1 Introduction

---

**FIGURE 100.27**  Root loci of Eq. (100.62).

## References and Further Information

R. C. Dorf, *Modern Control Systems,* 5th ed., Reading, Mass.: Addison-Wesley, 1989.

W. R. Evans, "Graphical analysis of control systems," *Trans. AIEE,* vol. 67, pp. 547–551, 1948.

B. C. Kuo, *Automatic Control Systems,* 6th ed., Englewood Cliffs, N.J.: Prentice-Hall, 1991a.

B. C. Kuo, *ACSP Software and Manual,* Englewood Cliffs, N.J.: Prentice-Hall, 1991b.

B. C. Kuo, *Digital Control Systems,* 2nd ed., New York: Holt, 1992a.

B. C. Kuo, *DCSP Software and Manual,* Champaign, Ill.: SRL, Inc., 1992b.

## 100.5   Compensation

*Charles L. Phillips and Royce D. Harbor*

**Compensation** is the process of modifying a closed-loop control system (usually by adding a *compensator* or *controller*) in such a way that the compensated system satisfies a given set of design specifications. This section presents the fundamentals of compensator design; actual techniques are available in the references.

A single-loop control system is shown in Fig. 100.28. This system has the transfer function from input $R(s)$ to output $C(s)$

$$T(s) = \frac{C(s)}{R(s)} = \frac{G_c(s)G_p(s)}{1 + G_c(s)G_p(s)H(s)} \qquad (100.63)$$

and the characteristic equation is

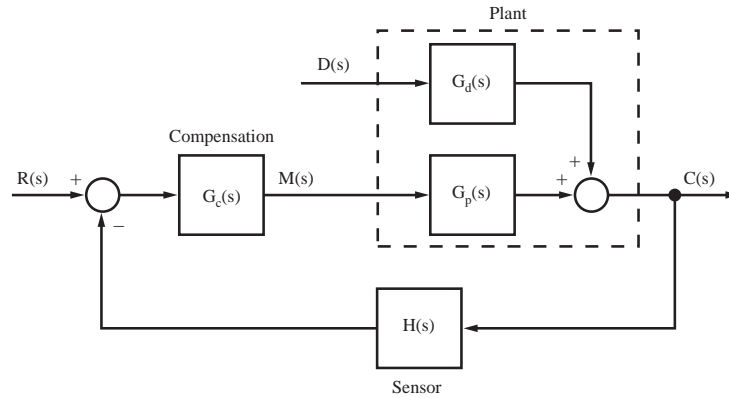$$1 + G_c(s)G_p(s)H(s) = 0 \qquad (100.64)$$

**FIGURE 100.28**   A closed-loop control system.

where $G_c(s)$ is the *compensator* transfer function, $G_p(s)$ is the *plant* transfer function, and $H(s)$ is the *sensor* transfer function. The transfer function from the disturbance input $D(s)$ to the output is $G_d(s)/[1 + G_c(s)G_p(s)H(s)]$. The function $G_c(s)G_p(s)H(s)$ is called the *open-loop function.*

## Control System Specifications

The compensator transfer function $G_c(s)$ is designed to give the closed-loop system certain specified characteristics, which are realized through achieving one or more of the following:

1. Improving the transient response. Increasing the speed of response is generally accomplished by increasing the open-loop gain $G_c(j\omega)G_p(j\omega)H(j\omega)$ at higher frequencies such that the system bandwidth is increased. Reducing overshoot (ringing) in the response generally involves increasing the phase margin $\phi_m$ of the system, which tends to remove any resonances in the system. The phase margin $\phi_m$ occurs at the frequency $\omega_1$ and is defined by the relationship

$$\left| G_c(j\omega_1)G_p(j\omega_1)H(j\omega_1) \right| = 1$$

   with the angle of $G_c(j\omega_1)G_p(j\omega_1)H(j\omega_1)$ equal to $(180° + \phi_m)$.
2. Reducing the steady-state errors. Steady-state errors are decreased by increasing the open-loop gain $G_c(j\omega)G_p(j\omega)H(j\omega)$ in the frequency range of the errors. Low-frequency errors are reduced by increasing the low-frequency open-loop gain and by increasing the type number of the system [the number of poles at the origin in the open-loop function $G_c(s)G_p(s)H(s)$].
3. Reducing the sensitivity to plant parameters. Increasing the open-loop gain $G_c(j\omega)G_p(j\omega)\ H(j\omega)$ tends to reduce the variations in the system characteristics due to variations in the parameters of the plant.
4. Rejecting disturbances. Increasing the open-loop gain $G_c(j\omega)G_p(j\omega)H(j\omega)$ tends to reduce the effects of disturbances [$D(s)$ in Fig. 100.28] on the system output, provided that the increase in gain does not appear in the direct path from disturbance inputs to the system output.
5. Increasing the relative stability. Increasing the open-loop gain tends to reduce phase and gain margins, which generally increases the overshoot in the system response. Hence, a trade-off exists between the beneficial effects of increasing the open-loop gain and the resulting detrimental effects of reducing the stability margins.

## Design

Design procedures for compensators are categorized as either *classical methods* or *modern methods.* Classical methods discussed are:

- Ph ase-lag frequency response
- Phase-lead frequency response
- Phase-lag root locus
- Phase-lead root locus

Modern methods discussed are:

- Pole placement
- State estimation
- Optimal

### Frequency Response Design

Classical design procedures are normally based on the open-loop function of the uncompensated system, $G_p(s)H(s)$. Two compensators are used in classical design; the first is called a *phase-lag compensator,* and the second is called a *phase-lead compensator.*

The general characteristics of phase-lag-compensated systems are as follows:

1. The low-frequency behavior of a system is improved. This improvement appears as reduced errors at low frequencies, improved rejection of low-frequency disturbances, and reduced sensitivity to plant parameters in the low-frequency region.
2. The system bandwidth is reduced, resulting in a slower system time response and better rejection of high-frequency noise in the sensor output signal.

The general characteristics of phase-lead-compensated systems are as follows:

1. The high-frequency behavior of a system is improved. This improvement appears as faster responses to inputs, improved rejection of high-frequency disturbances, and reduced sensitivity to changes in the plant parameters.
2. The system bandwidth is increased, which can increase the response to high-frequency noise in the sensor output signal.

The transfer function of a first-order compensator can be expressed as

$$G_c(s) = \frac{K_c(s/\omega_0 + 1)}{s/\omega_p + 1} \tag{100.65}$$

where $-\omega_0$ is the compensator zero, $-\omega_p$ is its pole, and $K_c$ is its dc gain. If $\omega_p < \omega_0$, the compensator is phase-lag. The Bode diagram of a phase-lag compensator is given in Fig. 100.29 for $K_c = 1$.

It is seen from Fig. 100.29 that the phase-lag compensator reduces the high-frequency gain of the open-loop function relative to the low-frequency gain. This effect allows a higher low-frequency gain, with the advantages listed above. The pole and zero of the compensator must be placed at very low frequencies relative to the compensated-system bandwidth so that the destabilizing effects of the negative phase of the compensator are negligible.

If $\omega_p > \omega_0$ the compensator is phase-lead. The Bode diagram of a phase-lead compensator is given in Fig. 100.30 for $K_c = 1$.

It is seen from Fig. 100.30 that the phase-lead compensator increases the high-frequency gain of the open-loop function relative to its low-frequency gain. Hence, the system has a larger bandwidth, with the advantages listed above. The pole and zero of the compensator are generally difficult to place, since the increased gain of the open-loop function tends to destabilize the system, while the phase lead of the compensator tends to stabilize the system. The pole-zero placement for the phase-lead compensator is much more critical than that of the phase-lag compensator.

A typical Nyquist diagram of an uncompensated system is given in Fig. 100.31. The pole and the zero of a phase-lag compensator are placed in the frequency band labeled *A*. This placement negates the destabilizing
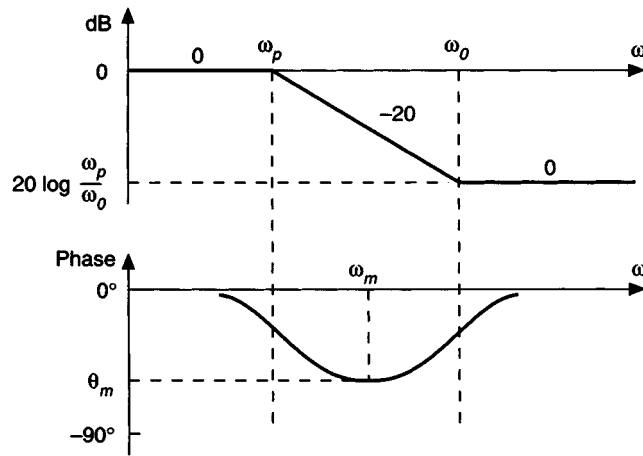
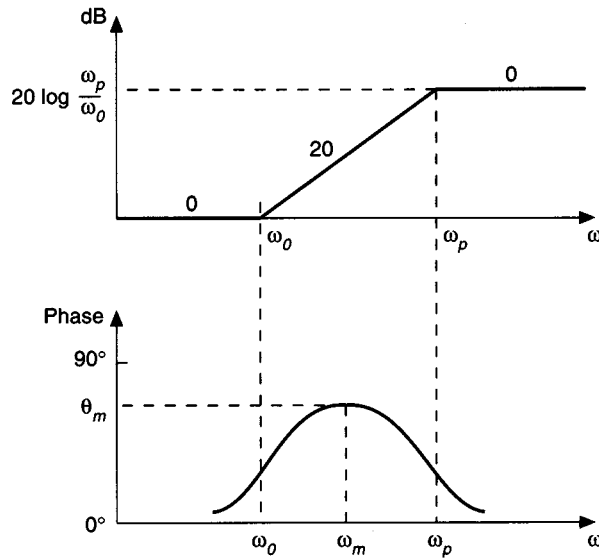**FIGURE 100.29**  Bode diagram for a phase-lag compensator.



**FIGURE 100.30**  Bode diagram for a phase-lead compensator.

effect of the negative phase of the compensator. The pole and zero of a phase-lead compensator are placed in the frequency band labeled *B*. This placement utilizes the stabilizing effect of the positive phase of the compensator.

### PID Controllers

Proportional-plus-integral-plus-derivative (PID) compensators are probably the most utilized form for compensators. These compensators are essentially equivalent to a phase-lag compensator cascaded with a phase-lead compensator. The transfer function of this compensator is given by

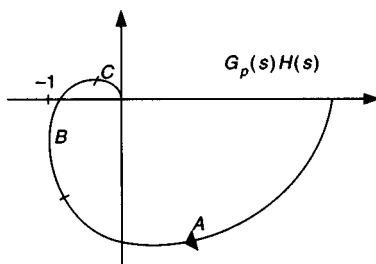$$G_c(s) = K_P + \frac{K_I}{s} + K_D s \qquad (100.66)$$

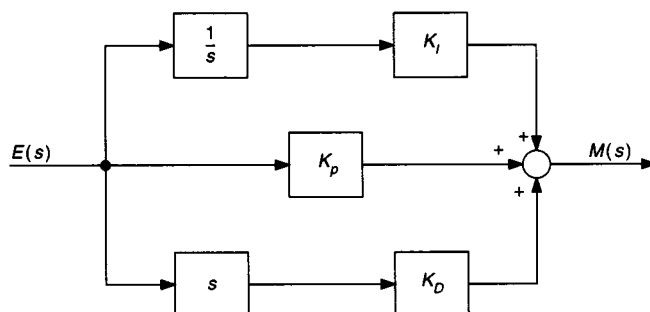**FIGURE 100.31**   A typical Nyquist diagram for $G_p(s)H(s)$.



**FIGURE 100.32**   Block diagram of a PID compensator.



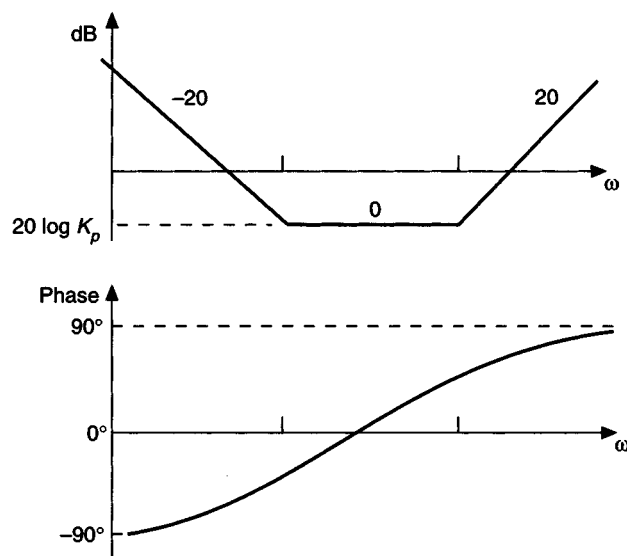**FIGURE 100.33**   Bode diagram of a PID compensator.

A block diagram portrayal of the compensator is shown in Fig. 100.32. The integrator in this compensator increases the system type by one, resulting in an improved low-frequency response. The Bode diagram of a PID compensator is given in Fig. 100.33.

With $K_D = 0$, the compensator is phase-lag, with the pole in (100.65) moved to $\omega_p = 0$. As a result the compensator is type one. The zero of the compensator is placed in the low-frequency range to correspond to the zero of the phase-lag compensator discussed above.

With $K_I = 0$, the compensator is phase-lead, with a single zero and the pole moved to infinity. Hence, the gain continues to increase with increasing frequency. If high-frequency noise is a problem, it may be necessary to add one or more poles to the PD or PID compensators. These poles must be placed at high frequencies relative to the phase-margin frequency such that the phase margin (stability characteristics) of the system is not degraded. PD compensators realized using rate sensors minimize noise problems [Phillips and Harbor, 1991].

### Root Locus Design

Root locus design procedures generally result in the placement of the two dominant poles of the closed-loop system transfer function. A system has two dominant poles if its behavior approximates that of a second-order system.

The differences in root locus designs and frequency response designs appear only in the interpretation of the control-system specifications. A root locus design that improves the low-frequency characteristics of the system will result in a phase-lag controller; a phase-lead compensator results if the design improves the high-frequency response of the system. If a root locus design is performed, the frequency response characteristics of the system should be investigated. Also, if a frequency response design is performed, the poles of the closed-loop transfer function should be calculated.

## Modern Control Design

The classical design procedures above are based on a transfer-function model of a system. Modern design procedures are based on a *state-variable model* of the plant. The plant transfer function is given by

$$\frac{Y(s)}{U(s)} = G_p(s) \tag{100.67}$$

where we use $u(t)$ for the plant input and $y(t)$ for the plant output. If the system model is $n$th order, the denominator of $G_p(s)$ is an $n$th-order polynomial.

The state-variable model, or state model, for a single-input–single-output plant is given by

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}u(t)$$
$$y(t) = \mathbf{C}\mathbf{x}(t) \tag{100.68}$$

$$y(t) = \mathbf{C}\mathbf{x}(t)$$

where $\mathbf{x}(t)$ is the $n \times 1$ state vector, $u(t)$ is the plant input, $y(t)$ is the plant output, $\mathbf{A}$ is the $n \times n$ *system matrix*, $\mathbf{B}$ is the $n \times 1$ *input matrix*, and $\mathbf{C}$ is the $1 \times n$ *output matrix*. The transfer function of (100.67) is an input-output model; the state model of (100.68) yields the same input-output model and in addition includes an internal model of the system. The state model of (100.68) is readily adaptable to a multiple-input–multiple-output system (*a multivariable system*); for that case, $u(t)$ and $y(t)$ are vectors. We will consider only single-input–single-output systems.

The plant transfer function of (100.67) is related to the state model of (100.68) by

$$G_p(s) = \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} \tag{100.69}$$

The state model is not unique; many combinations of the matrices $\mathbf{A}$, $\mathbf{B}$, and $\mathbf{C}$ can be found to satisfy (100.69) for a given transfer function $G_p(s)$.
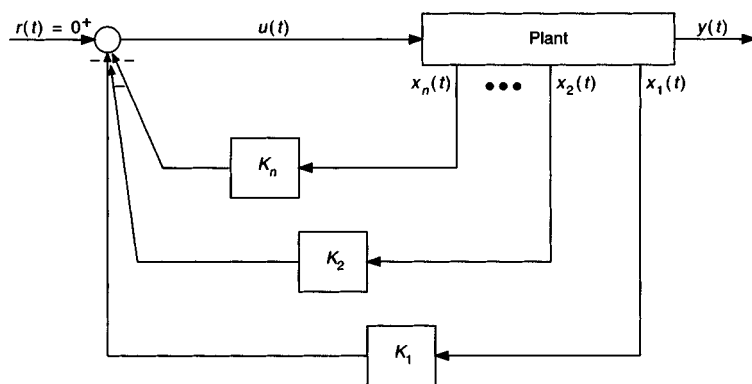
**FIGURE 100.34**  Implementation of pole-placement design.

Classical compensator design procedures are based on the open-loop function $G_p(s)H(s)$ of Fig. 100.28. It is common to present modern design procedures as being based on only the plant model of (100.68). However, the models of the sensors that measure the signals for feedback must be included in the state model. This problem will become more evident as the modern procedures are presented.

### Pole Placement

Probably the simplest modern design procedure is *pole placement*. Recall that root locus design was presented as placing the two dominant poles of the closed-loop transfer function at desired locations. The pole-placement procedure places *all* poles of the closed-loop transfer function, or equivalently, all roots of the closed-loop system characteristic equation, at desirable locations.

The system design specifications are used to generate the desired closed-loop system characteristic equation $\alpha_c(s)$, where

$$\alpha_c(s) = s^n + \alpha_{n-1}s^{n-1} + \ldots + \alpha_1 s + \alpha_0 = 0 \qquad (100.70)$$

for an $n$th-order plant. This characteristic equation is realized by requiring the plant input to be a linear combination of the plant states, that is,

$$u(t) = -K_1 x_1(t) - K_2 x_2(t) - \ldots - K_n x_n(t) = -\mathbf{K}\mathbf{x}(t) \qquad (100.71)$$

where $\mathbf{K}$ is the $1 \times n$ feedback-gain matrix. Hence *all* states must be measured and fed back. This operation is depicted in Fig. 100.34.

The feedback-gain matrix $\mathbf{K}$ is determined from the desired characteristic equation for the closed-loop system of (100.70):

$$\alpha_c(s) = \left| sI - A + BK \right| = 0 \qquad (100.72)$$

The state feedback gain matrix $\mathbf{K}$ which yields the specified closed-loop characteristic equation $\alpha_c(s)$ is

$$\mathbf{K} = [0\ 0\ \ldots\ 0\ 1][\mathbf{B}\ \mathbf{AB}\ \ldots\ \mathbf{A}^{n-1}\mathbf{B}]^{-1}\alpha_c(\mathbf{A}) \qquad (100.73)$$

where $\alpha_c(\mathbf{A})$ is (100.70) with the scalar $s$ replaced with the matrix $\mathbf{A}$. A plant is said to be *controllable* if the inverse matrix in (100.73) exists. Calculation of $\mathbf{K}$ completes the design process. A simple computer algorithm is available for solving (100.73) for $\mathbf{K}$.
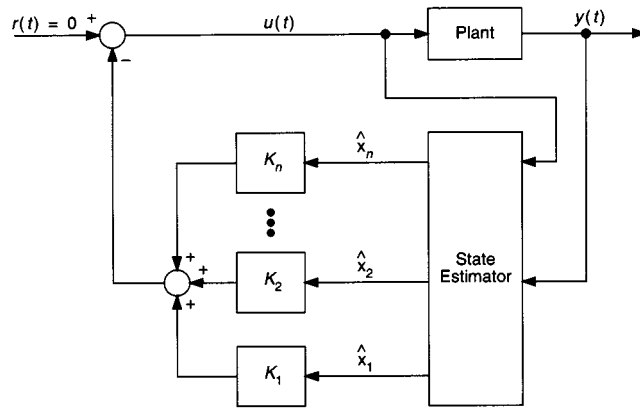
**FIGURE 100.35** Implementation of observer-pole-placement design.

## State Estimation

In general, modern design procedures require that the state vector $\mathbf{x}(t)$ be fed back, as in (100.71). The measurement of all state variables is difficult to implement for high-order systems. The usual procedure is to estimate the states of the system from the measurement of the output $y(t)$, with the estimated states then fed back.

Let the estimated state vector be $\hat{\mathbf{x}}$. One procedure for estimating the system states is by an *observer,* which is a dynamic system realized by the equations

$$\frac{d\hat{\mathbf{x}}(t)}{dt} = (\mathbf{A} - \mathbf{GC})\hat{\mathbf{x}}(t) + \mathbf{B}u(t) + \mathbf{G}y(t) \qquad (100.74)$$

with the feedback equation of (100.71) now realized by

$$u(t) = -\mathbf{K}\hat{\mathbf{x}}(t) \qquad (100.75)$$

The matrix $\mathbf{G}$ in (100.74) is calculated by assuming an $n$th-order characteristic equation for the observer of the form

$$\alpha_e(s) = \left| sI - A + GC \right| = 0 \qquad (100.76)$$

The estimator gain matrix $\mathbf{G}$ which yields the specified estimator characteristic equation $\alpha_e(s)$ is

$$\mathbf{G} = \alpha_e(\mathbf{A})[\mathbf{C} \ \mathbf{CA} \ \ldots \ \mathbf{CA}^{n-1}]^{-T}[0 \ 0 \ \ldots \ 0 \ 1]^T \qquad (100.77)$$

where $[\cdot]^T$ denotes the matrix transpose. A plant is said to be *observable* if the inverse matrix in (100.77) exists. An implementation of the closed-loop system is shown in Fig. 100.35. The observer is usually implemented on a digital computer. The plant and the observer in Fig. 100.35 are both $n$th-order; hence, the closed-loop system is of order $2n$.

The observer-pole-placement system of Fig. 100.35 is equivalent to the system of Fig. 100.36, which is of the form of closed-loop systems designed by classical procedures. The transfer function of the controller-estimator (equivalent compensator) of Fig. 100.36 is given by

$$\mathbf{G}_{ec}(s) = \mathbf{K}[s\mathbf{I} - \mathbf{A} + \mathbf{GC} + \mathbf{BK}]^{-1}\mathbf{G} \qquad (100.78)$$
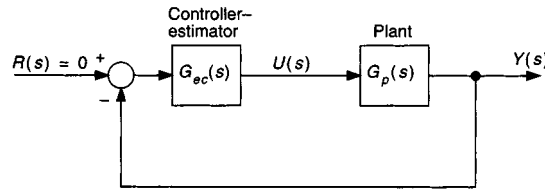
**FIGURE 100.36** Equivalent system for pole-placement design.

This compensator is $n$th-order for an $n$th-order plant; hence, the total system is of order $2n$. The characteristic equation for the compensated system is given by

$$|s\mathbf{I} - \mathbf{A} + \mathbf{BK}| \, |s\mathbf{I} - \mathbf{A} + \mathbf{GC}| = \alpha_c(s)\alpha_e(s) = 0 \tag{100.79}$$

The roots of this equation are the roots of the pole-placement design plus those of the observer design. For this reason, the roots of the characteristic equation for the observer are usually chosen to be faster than those of the pole-placement design.

### Linear Quadratic Optimal Control

We define an optimal control system as one for which some mathematical function is minimized. The function to be minimized is called the *cost function*. For steady-state linear quadratic optimal control the cost function is given by

$$V_\infty = \int_t^\infty [\mathbf{x}^T(\tau)\mathbf{Q}\mathbf{x}(\tau) + Ru^2(\tau)]\, d\tau \tag{100.80}$$

where Q and R are chosen to satisfy the design criteria. In general, the choices are not straightforward. Minimization of (100.80) requires that the plant input be given by

$$u(t) = -R^{-1}\mathbf{B}^T\mathbf{M}_\infty\mathbf{x}(t) \tag{100.81}$$

where the $n \times n$ matrix $\mathbf{M}_\infty$ is the solution to the *algebraic Riccati equation*

$$\mathbf{M}_\infty\mathbf{A} + \mathbf{A}^T\mathbf{M}_\infty - \mathbf{M}_\infty\mathbf{BR}^{-1}\mathbf{B}^T\mathbf{M}_\infty + \mathbf{Q} = 0 \tag{100.82}$$

The existence of a solution for this equation is involved [Friedland, 1986] and is not presented here. Optimal control systems can be designed for cost functions other than that of (100.80).

## Other Modern Design Procedures

Other modern design procedures exist; for example, *self-tuning control systems* continually estimate certain plant parameters and adjust the compensator based on this estimation. These control systems are a type of *adaptive control systems* and usually require that the control algorithms be implemented using a digital computer. These control systems are beyond the scope of this book (see, for example, Astrom and Wittenmark, 1984).

## Defining Term

**Compensation:**   The process of physically altering a closed-loop system such that the system has specified characteristics. This alteration is achieved either by changing certain parameters in the system or by adding a physical system to the closed-loop system; in some cases both methods are used.

## Related Topics

100.3 Frequency Response Methods: Bode Diagram Approach • 100.4 Root Locus

## References

K. J. Astrom and B. Wittenmark, *Computer Controlled Systems,* Englewood Cliffs, N.J.: Prentice-Hall, 1984.

W. L. Brogan, *Modern Control Theory,* Englewood Cliffs, N.J.: Prentice-Hall, 1985.

R. C. Dorf, *Modern Control Systems,* 7th ed., Reading, Mass.: Addison-Wesley, 1995.

G. F. Franklin, J. D. Powell, and A. Emami-Naeini, *Feedback Control of Dynamic Systems,* Reading, Mass.: Addison-Wesley, 1986.

B. Friedland, *Control System Design,* New York: McGraw-Hill, 1986.

B. C. Kuo, *Automatic Control Systems,* Englewood Cliffs, N.J.: Prentice-Hall, 1987.

C. L. Phillips and R. D. Harbor, *Feedback Control Systems,* 2nd ed., Englewood Cliffs, N.J.: Prentice-Hall, 1991.

## 100.6 Digital Control Systems

*Raymond G. Jacquot and John E. McInroy*

The use of the **digital computer** to control physical processes has been a topic of discussion in the technical literature for over four decades, but the actual use of a digital computer for control of industrial processes was reserved only for massive and slowly varying processes such that the high cost and slow computing speed of available computers could be tolerated. The invention of the integrated circuit microprocessor in the early 1970s radically changed all that; now microprocessors are used in control tasks in automobiles and household appliances, applications where high cost is not justifiable.

When the term *digital control* is used, it usually refers to the process of employing a digital computer to control some process that is characterized by continuous-in-time dynamics. The control can be of the open-loop variety where the control strategy output by the digital computer is dictated without regard to the status of the process variables. An alternative technique is to supply the digital computer with digital data about the process variables to be controlled, and thus the control strategy output by the computer depends on the process variables that are to be controlled. This latter strategy is a **feedback control** strategy wherein the computer, the process, and interface hardware form a closed loop of information flow.

Examples of dynamic systems that are controlled in such a closed-loop digital fashion are flight control of civilian and military aircraft, control of process variables in chemical processing plants, and position and force control in industrial robot manipulators. The simplest form of feedback control strategy provides an on-off control to the controlling variables based on measured values of the process variables. This strategy will be illustrated by a simple example in a following subsection.

In the past decade and a half many excellent textbooks on the subject of digital control systems have been written, and most of them are in their second edition. The texts in the References provide in-depth development of the theory by which such systems are analyzed and designed.

### A Simple Example

Such a closed-loop or feedback control situation is illustrated in Fig. 100.37, which illustrates the feedback control of the temperature in a simple environmental chamber that is to be kept at a constant temperature somewhat above room temperature.

Heat is provided by turning on a relay that supplies power to a heater coil. The on-off signal to the relay can be supplied by 1 bit of an output port of the microprocessor (typically the port would be 8 bits wide). A second bit of the port can be used to turn a fan on and off to supply cooling air to the chamber. An analog-to-digital (A/D) converter is employed to convert the amplified thermocouple signal to a digital word that is then supplied to the input port of the microprocessor. The program being executed by the microprocessor reads the temperature data supplied to the input port and compares the binary number representing the temperature to a binary
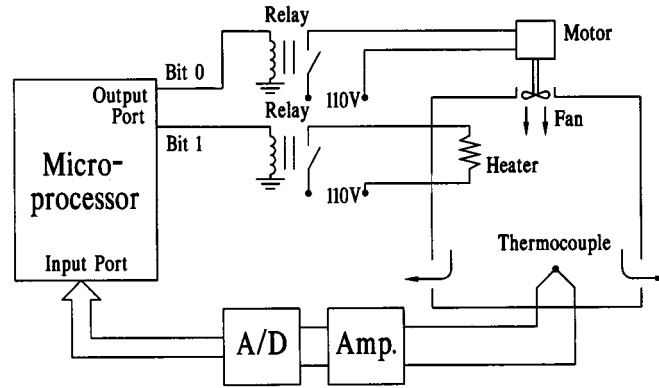
**FIGURE 100.37**  Microprocessor control of temperature in a simple environmental chamber.
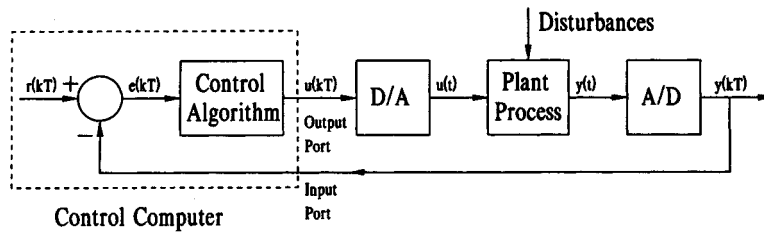


**FIGURE 100.38**  Closed-loop control of a single process variable.

version of the desired temperature and makes a decision whether or not to turn on the heater or the fan or to do nothing. The program being executed runs in a continuous loop, repeating the operations discussed above.

This simple on-off control strategy is often not the best when extremely precise control of the process variables is required. A more precise control may be obtained if the controlling variable levels can be adjusted to be somewhat larger if the deviation of the process variable from the desired value is larger.

## Single-Loop Linear Control Laws

Consider the case where a single variable of the process is to be controlled, as illustrated in Fig. 100.38. The output of the plant $y(t)$ is to be sampled every $T$ seconds by an A/D converter, and this sequence of numbers will be denoted as $y(kT)$, $k = 0, 1, 2, \ldots$. The goal is to make the sequence $y(kT)$ follow some desired known sequence [the reference sequence $r(kT)$]. Consequently, the sequence $y(kT)$ is subtracted from $r(kT)$ to obtain the so-called error sequence $e(kT)$. The control computer then acts on the error sequence, using some control algorithms, to produce the control effort sequence $u(kT)$ that is supplied to the digital-to-analog (D/A) converter which then drives the actuating hardware with a signal proportional to $u(kT)$. The output of the D/A converter is then held constant on the current time interval, and the control computer waits for the next sample of the variable to be controlled, the arrival of which repeats the sequence. The most commonly employed control algorithm or control law is a linear difference equation of the form

$$u(kT) = a_n e(kT) + a_{n-1} e((k-1)T) + \ldots + a_0 e((k-n)T)$$

$$+ b_{n-1} u((k-1)T) + \ldots + b_0 u((k-n)T) \tag{100.83}$$

The question remains as to how to select the coefficients $a_0, \ldots, a_n$ and $b_0, \ldots, b_{n-1}$ in expression (100.83) to give an acceptable degree of control of the plant.

## Proportional Control

This is the simplest possible control algorithm for the digital processor wherein the most current control effort is proportional to the current error or using only the first term of relation (100.83)

$$u(kT) = a_n e(kT) \tag{100.84}$$

This algorithm has the advantage that it is simple to program, while, on the other hand, its disadvantage lies in the fact that it has poor disturbance rejection properties in that if $a_n$ is made large enough for good disturbance rejection, the closed-loop system can be unstable (i.e., have transient responses which increase with time). Since the object is to regulate the system output in a known way, these unbounded responses preclude this regulation.

## PID Control Algorithm

A common technique employed for decades in chemical process control loops is that of proportional-plus-integral-plus-derivative (PID) control wherein a continuous-time control law would be given by

$$u(t) = K_p e(t) + K_i \int_0^t e(\tau)d\tau + K_d \frac{de}{dt} \tag{100.85}$$

This would have to be implemented by an analog filter.

To implement the design in digital form the proportional term can be carried forward as in relation (100.84); however, the integral can be replaced by trapezoidal integration using the error sequence, while the derivative can be replaced with the backward difference resulting in a computer control law of the form [Jacquot, 1995]

$$u(kT) = u((k-1)T) + \left( K_p + \frac{K_i T}{2} + \frac{K_d}{T} \right) e(kT)$$
$$+ \left( \frac{K_i T}{2} - K_p - \frac{2K_d}{T} \right) e((k-1)T) + \frac{K_d}{T} e((k-2)T) \tag{100.86}$$

where $T$ is the duration of the sampling interval. The selection of the coefficients in this algorithm ($K_i$, $K_d$, and $K_p$) is best accomplished by the Ziegler-Nichols tuning process [Franklin et al., 1990].

## The Closed-Loop System

When the plant process is linear or may be linearized about an operating point and the control law is linear as in expressions (100.83), (100.84), or (100.86), then an appropriate representation of the complete closed-loop system is by the so-called $z$-transform. The $z$-transform plays the role for linear, constant-coefficient difference equations that the Laplace transform plays for linear, constant-coefficient differential equations. This $z$-domain representation allows the system designer to investigate system time response, frequency response, and stability in a single analytical framework.

If the plant can be represented by an $s$-domain transfer function $G(s)$, then the discrete-time ($z$-domain) transfer function of the plant, the analog-to-digital converter, and the driving digital-to-analog converter is

$$G(z) = \left( \frac{z-1}{z} \right) Z \left\{ L^{-1} \left[ \frac{G(s)}{s} \right] \right\} \tag{100.87}$$
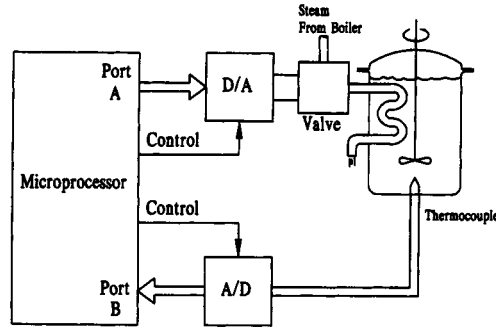
**FIGURE 100.39**  A computer-controlled thermal mixing tank.

where $Z(\cdot)$ is the $z$-transform and $L^{-1}(\cdot)$ is the inverse Laplace transform. The transfer function of the control law of (100.83) is

$$D(z) = \frac{U(z)}{E(z)} = \frac{a_n z^n + a_{n-1} z^{n-1} + \cdots + a_0}{z^n - b_{n-1} z^{n-1} - \cdots - b_0} \tag{100.88}$$

For the closed-loop system of Fig. 100.38 the closed-loop $z$-domain transfer function is

$$M(z) = \frac{Y(z)}{R(z)} = \frac{G(z)D(z)}{1 + G(z)D(z)} \tag{100.89}$$

where $G(z)$ and $D(z)$ are specified above. The characteristic equation of the closed-loop system is

$$1 + G(z)D(z) = 0 \tag{100.90}$$

The dynamics and stability of the system can be assessed by the locations of the zeros of (100.90) (the closed-loop poles) in the complex $z$ plane. For stability the zeros of (100.90) above must be restricted to the unit circle of the complex $z$ plane.

## A Linear Control Example

Consider the temperature control of a chemical mixing tank shown in Fig. 100.39. From a transient power balance the differential equation relating the rate of heat added $q(t)$ to the deviation in temperature from the ambient $\theta(t)$ is given as

$$\frac{d\theta}{dt} + \frac{1}{\tau}\theta = \frac{1}{mc}q(t) \tag{100.91}$$

where $\tau$ is the time constant of the process and $mc$ is the heat capacity of the tank. The transfer function of the tank is

$$\frac{\Theta(s)}{Q(s)} = G(s) = \frac{1/mc}{s + 1/\tau} \tag{100.92}$$

The heater is driven by a D/A converter, and the temperature measurement is sampled with an A/D converter. The data converters are assumed to operate synchronously, so the discrete-time transfer function of the tank and the two data converters is from expression (100.87):

$$G(z) = \frac{\Theta(z)}{Q(z)} = \frac{\tau}{mc} \frac{1 - e^{-T/\tau}}{z - e^{-T/\tau}} \qquad (100.93)$$

If a proportional control law is chosen, the transfer function associated with the control law is the gain $a_n = K$ or

$$D(z) = K \qquad (100.94)$$

The closed-loop characteristic equation is from (100.90):

$$1 + \frac{K\tau}{mc} \frac{1 - e^{-T/\tau}}{z - e^{-T/\tau}} = 0 \qquad (100.95)$$

If a common denominator is found, the resulting numerator is

$$z - e^{-T/\tau} + \frac{K\tau}{mc}(1 - e^{-T/\tau}) = 0 \qquad (100.96)$$

The root of this equation is

$$z = e^{-T/\tau} + \frac{K\tau}{mc}(e^{-T/\tau} - 1) \qquad (100.97)$$

If this root location is investigated as the gain parameter $K$ is varied upward from zero, it is seen that the root starts at $z = e^{-T/\tau}$ for $K = 0$ and moves to the left along the real axis as $K$ increases. Initially it is seen that the system becomes faster, but at some point the responses become damped and oscillatory, and as $K$ is further increased the oscillatory tendency becomes less damped, and finally a value of $K$ is reached where the oscillations are sustained at constant amplitude. A further increase in $K$ will yield oscillations that increase with time. Typical unit step responses for $r(k) = 1$ and $T/\tau = 0.2$ are shown in Fig. 100.40.

It is easy to observe this tendency toward oscillation as $K$ increases, but a problem that is clear from Fig. 100.40 is that in the steady state there is a persistent error between the response and the reference [$r(k) = 1$]. Increasing the gain $K$ will make this error smaller at the expense of more oscillations. As a remedy for this steady-state error problem and control of the dynamics, a control law transfer function $D(z)$ will be sought that inserts integrator action into the loop while simultaneously canceling the pole of the plant. This dictates that the controller have a transfer function of the form

$$D(z) = \frac{U(z)}{E(z)} = \frac{K(z - e^{-T/\tau})}{z - 1} \qquad (100.98)$$

Typical unit step responses are illustrated in Fig. 100.41 for several values of the gain parameter. The control law that must be programmed in the digital processor is

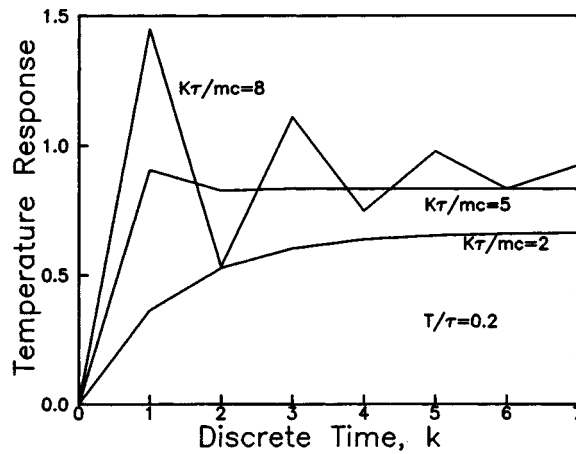$$u(kT) = u((k-1)T) + K[e(kT) - e^{-T/\tau}e((k-1)T)] \qquad (100.99)$$

**FIGURE 100.40**   Step responses of proportionally controlled thermal mixing tank.
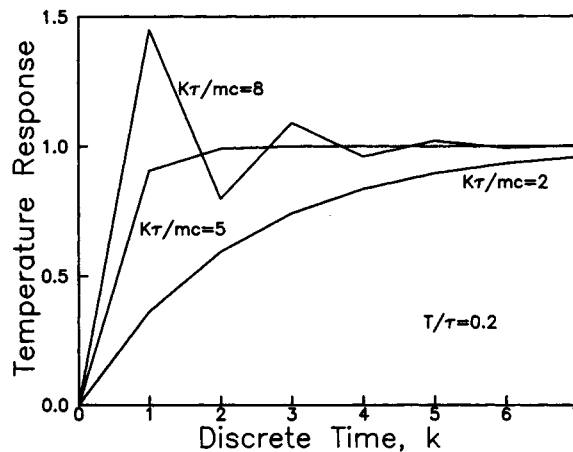


**FIGURE 100.41**   Step responses of the compensated thermal mixing tank.

The additional effort to program this over that required to program the proportional control law of (100.94) is easily justified since $K$ and $e^{-T/\tau}$ are simply constants.

## Defining Terms

**Digital computer:**   A collection of digital devices including an arithmetic logic unit (ALU), read-only memory (ROM), random-access memory (RAM), and control and interface hardware.

**Feedback control:**   The regulation of a response variable of a system in a desired manner using measurements of that variable in the generation of the strategy of manipulation of the controlling variables.

## Related Topics

8.1 Introduction • 112.1 Introduction • 112.3 The State of the Art in CACSD

## References

K.J. Astrom and B. Wittenmark, *Computer Controlled Systems: Theory and Design,* Englewood Cliffs, N.J.: Prentice-Hall, 1984.

G.F. Franklin, J.D. Powell, and M.L. Workman, *Digital Control of Dynamic Systems,* 2nd ed., Reading, Mass.: Addison-Wesley, 1990.

C.H. Houpis and G.B. Lamont, *Digital Control Systems: Theory, Hardware, Software,* 2nd ed., New York: McGraw-Hill, 1992.

R.G. Jacquot, *Modern Digital Control Systems,* 2nd ed., New York: Marcel Dekker, 1995.

B.C. Kuo, *Digital Control Systems,* 2nd ed., Orlando, Fla.: Saunders, 1992.

C.L. Phillips and H.T. Nagle, *Digital Control System Analysis and Design,* 3rd ed., Englewood Cliffs, N.J.: Prentice-Hall, 1995.

R. J. Vaccaro, *Digital: A State-Space Approach,* New York: McGraw-Hill, 1995.

## Further Information

The *IEEE Control Systems Magazine* is a useful information source on control systems in general and digital control in particular. Highly technical articles on the state of the art in digital control may be found in the *IEEE Transactions on Automatic Control, the IEEE Transactions on Control Systems Technology,* and the ASME *Journal of Dynamic Systems, Measurement and Control.*

## 100.7   Nonlinear Control Systems[3]

*Derek P. Atherton*

### The Describing Function Method

The describing function method, abbreviated as DF, was developed in several countries in the 1940s [Atherton, 1982], to answer the question: "What are the necessary and sufficient conditions for the nonlinear feedback system of Fig. 100.42 to be stable?" The problem still remains unanswered for a system with static nonlinearity, $n(x)$, and linear plant $G(s)$. All of the original investigators found limit cycles in control systems and observed that, in many instances with structures such as Fig. 100.42, the wave form of the oscillation at the input to the nonlinearity was almost sinusoidal. If, for example, the nonlinearity in Fig. 100.42 is an ideal relay, that is has an on-off characteristic, so that an odd symmetrical input wave form will produce a square wave at its output, the output of $G(s)$ will be almost sinusoidal when $G(s)$ is a low pass filter which attenuates the higher harmonics in the square wave much more than the fundamental. It was, therefore, proposed that the nonlinearity should be represented by its gain to a sinusoid and that the conditions for sustaining a sinusoidal limit cycle be evaluated to assess the stability of the feedback loop. Because of the nonlinearity, this gain in response to a sinusoid is a function of the amplitude of the sinusoid and is known as the describing function. Because describing function methods can be used other than for a single sinusoidal input, the technique is referred to as the single sinusoidal DF or sinusoidal DF.

### The Sinusoidal Describing Function

For the reasons explained above, if we assume in Fig. 100.42 that $x(t) = a \cos \theta$, where $\theta = \omega t$ and $n(x)$ is a symmetrical odd nonlinearity, then the output $y(t)$ will be given by the Fourier series,

$$y(\theta) = \sum_{n=0}^{\infty} a_n \cos n\theta + b_n \sin n\theta, \qquad (100.100)$$

$$\text{where } a_0 = 0, \qquad (100.101)$$

---

[3]The material in this section was previously published by CRC Press in *The Control Handbook,* William S. Levine, Ed., 1996.
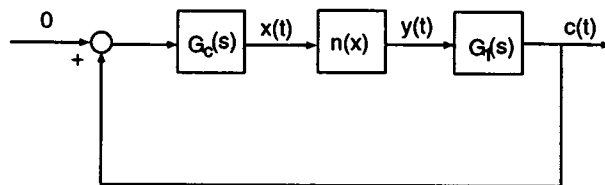
**FIGURE 100.42** Block diagram of a nonlinear system.

$$a_1 = (1/\pi)\int_0^{2\pi} y(\theta)\cos\theta d\theta, \qquad (100.102)$$

and

$$b_1 = (1/\pi)\int_0^{2\pi} y(\theta)\sin\theta d\theta. \qquad (100.103)$$

The fundamental output from the nonlinearity is $a_1\cos\theta + b_1\sin\theta$, so that the describing function, DF, defined as the fundamental output divided by the input amplitude, is complex and given by

$$N(a) = (a_1 - jb_1)/a \qquad (100.104)$$

which may be written

$$N(a) = N_p(a) + jN_q(a) \qquad (100.105)$$

where

$$N_p(a) = a_1/a \text{ and } N_q(a) = -b_1/a. \qquad (100.106)$$

Alternatively, in polar coordinates,

$$N(a) = M(a)e^{j\psi(a)} \qquad (100.107)$$

where

$$M(a) = (a_1^2 + b_1^2)^{1/2}/a$$

and

$$\Psi(a) = -\tan^{-1}(b_1/a_1). \qquad (100.108)$$

If $n(x)$ is single valued, then $b_1 = 0$ and

$$a_1 = (4/\pi)\int_0^{\pi/2} y(\theta)\cos\theta d\theta \qquad (100.109)$$

giving

$$N(a) = a_1/a = (4/a\pi)\int_0^{\pi/2} y(\theta) \cos \theta d\theta \qquad (100.110)$$

Although Eqs. (100.102) and (100.103) are an obvious approach to evaluating the fundamental output of a nonlinearity, they are indirect, because one must first determine the output wave form $y(\theta)$ from the known nonlinear characteristic and sinusoidal input wave form. This is avoided if the substitution $\theta = \cos^{-1}(x/a)$ is made. After some simple manipulations,

$$a_1 = (4/a)\int_0^a x n_p(x) p(x) dx \qquad (100.111)$$

and

$$b_1 = (4/a\pi)\int_0^a n_q(x) dx. \qquad (100.112)$$

The function $p(x)$ is the amplitude probability density function of the input sinusoidal signal given by

$$p(x) = (1/\pi)(a^2 - x^2)^{-1/2}. \qquad (100.113)$$

The nonlinear characteristics $n_p(x)$ and $n_q(x)$, called the inphase and quadrature nonlinearities, are defined by

$$n_p(x) = [n_1(x) + n_2(x)]/2 \qquad (100.114)$$

and

$$n_q(x) = [n_2(x) - n_1(x)]/2 \qquad (100.115)$$

where $n_1(x)$ and $n_2(x)$ are the portions of a double-valued characteristic traversed by the input for $\dot{x} > 0$ and $\dot{x} < 0$, respectively, When the nonlinear characteristic is single-valued, $n_1(x) = n_2(x)$, so $n_p(x) = n(x)$ and $n_q(x) = 0$. Integrating Eq. (100.111) by parts yields

$$a_1 = (4/\pi)n(0^+) + (4/a\pi)\int_0^a n'(x)(a^2 - x^2)^{1/2} dx \qquad (100.116)$$

where $n'(x) = dn(x)/dx$ and $n(0^+) = \lim_{\epsilon \to \infty} n(\epsilon)$, a useful alternative expression for evaluating $a_1$.

An additional advantage of using Eqs. (100.111) and (100.112) is that they yield proofs of some properties of the DF for symmetrical odd nonlinearities. These include the following:

1. For a double-valued nonlinearity, the quadrature component $N_q(a)$ is proportional to the area of the nonlinearity loop, that is,

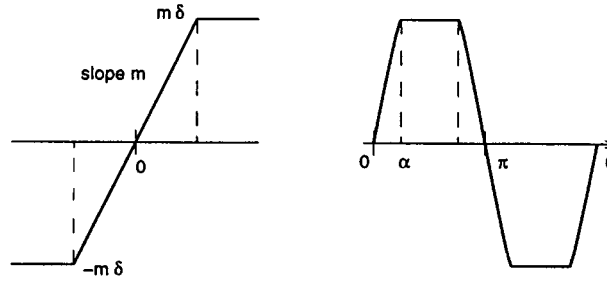$$N_q(a) = -(1/a^2\pi)(\text{area of nonlinearity loop}) \qquad (100.117)$$

**FIGURE 100.43**  Saturation nonlinearity.

2. For two single-valued nonlinearities $n_\alpha(x)$ and $n_\beta(x)$, with $n_\alpha(x) < n_\beta(x)$ for all $0 < x < b$, $N_\alpha(a) < N_\beta(a)$ for input amplitudes less than $b$.
3. For a single-valued nonlinearity with $k_1 x < n(x) < k_2 x$ for all $0 < x < b$, $k_1 < N(a) < k_2$ for input amplitudes less than $b$. This is the sector property of the DF; a similar result can be obtained for a double-valued nonlinearity [Cook, 1973].

When the nonlinearity is single valued, from the properties of Fourier series, the DF, $N(a)$, may also be defined as:

1. the variable gain, $K$, having the same sinusoidal input as the nonlinearity, which minimizes the mean squared value of the error between the output from the nonlinearity and that from the variable gain, and
2. the covariance of the input sinusoid and the nonlinearity output divided by the variance of the input.

## Evaluation of the Describing Function

To illustrate the evaluation of the DF two simple examples are considered.

### Saturation Nonlinearity

To calculate the DF, the input can alternatively be taken as $a \sin \theta$. For an ideal saturation characteristic, the nonlinearity output wave form $y(\theta)$ is as shown in Fig. 100.43. Because of the symmetry of the nonlinearity, the fundamental of the output can be evaluated from the integral over a quarter period so that

$$N(a) = \frac{4}{a\pi} \int_0^{\pi/2} y(\theta) \sin \theta d\theta,$$

which, for $a > \delta$, gives

$$N(a) = \frac{4}{a\pi} \left[ \int_0^\alpha ma \sin^2 \theta d\theta + \int_\alpha^{\pi/2} m\delta \sin \theta d\theta \right]$$

where $a = \sin^{-1} \delta/a$. Evaluation of the integrals gives

$$N(a) = (4m/\pi) \left[ \frac{\alpha}{2} - \frac{\sin 2\alpha}{4} + \delta \cos \alpha \right]$$

which, on substituting for $\delta$, give the result

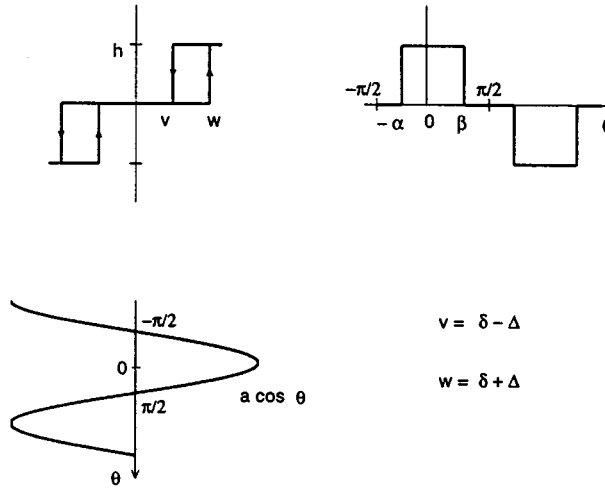$$N(a) = (m/\pi)(2\alpha + \sin 2\alpha). \tag{100.118}$$

**FIGURE 100.44**   Relay with dead zone and hysteresis.

Because, for $a < \delta$, the characteristic is linear giving $N(a) = m$, the DF for ideal saturation is $mN_s(\delta/a)$ where

$$N_s(\delta/a) = \begin{cases} 1, & \text{for } a < \delta, \text{ and} \\ (1/\pi)[2\alpha + \sin 2\alpha], & \text{for } a > \delta, \end{cases} \qquad (100.119)$$

where $a = \sin^{-1} \delta/a$.

Alternatively, one can evaluate $N(a)$ from Eq. (100.116), yielding

$$N(a) = a_1/a = (4/a^2\ \pi)\int_0^\delta m(a^2 - x^2)^{1/2}\,dx.$$

Using the substitution $x = a \sin \theta$,

$$N(a) = (4m/\pi)\int_0^\alpha \cos^2 \theta\, d\theta = (m/\pi)(2\alpha + \sin 2\alpha)$$

as before.

### Relay with Dead Zone and Hysteresis

The characteristic is shown in Fig. 100.44 together with the corresponding input, assumed equal to $a \cos \theta$, and the corresponding output wave form. Using Eqs. (100.102) and (100.103) over the interval $-\pi/2$ to $\pi/2$ and assuming that the input amplitude $a$ is greater than $\delta + \Delta$,

$$a_1 = \left(2/\pi \int_{-\alpha}^\beta h \cos \theta\, d\theta\right)$$

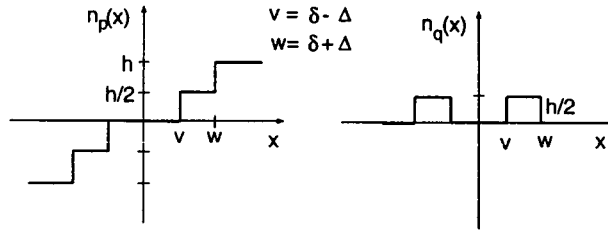$$= (2h/\pi)(\sin \beta + \sin \alpha),$$

**FIGURE 100.45** Function $n_p(x)$ and $n_q(x)$ for the relay of Figure 100.44.

where $\alpha = \cos^{-1}\left[(\delta - \Delta)/a\right]$ and $\beta = \cos^{-1}\left[(\delta + \Delta)/a\right]$, and

$$b_1 = \left(2/\pi\right)\int_{-\alpha}^{\beta} h\,\sin\theta\,d\theta$$

$$= \left(-2h/\pi\right)\left(\frac{(\delta + \Delta)}{a} - \frac{\delta - \Delta}{a}\right) = 4h\Delta/a\pi\ .$$

Thus

$$N(a) = \frac{2h}{a^2\pi}\left\{\left[a^2 - (\delta + \Delta)^2\right]^{1/2} + \left[a^2 - (\delta - \Delta)^2\right]^{1/2}\right\} - \frac{j4h\Delta}{a^2\pi}. \qquad (100.120)$$

For the alternative approach, one must first obtain the in-phase and quadrature nonlinearities shown in Fig. 100.45. Using Eqs. (100.111) and (100.112),

$$a_1 = \left(4/a\right)\int_{\delta-\Delta}^{\delta+\Delta} x\left(h/2\right)p(x)dx + \int_{\delta+\Delta}^{a} xhp(x)dx,$$

$$= \frac{2h}{a\pi}\left\{\left[a^2 - (\delta + \Delta)^2\right]^{1/2} + \left[a^2 - (\delta - \Delta)^2\right]^{1/2}\right\},$$

and

$$b_1 = \left(4/a\pi\right)\int_{d-\Delta}^{\delta+\Delta} \left(h/2\ dx\right) = 4h\Delta/a\pi$$

$$= \left(\text{Area of nonlinearity loop}\right)/a\pi$$

as before.

The DF of two nonlinearities in parallel equals the sum of their individual DFs, a result very useful for determining DFs, particularly of linear segmented characteristics with multiple break points. Several procedures [Altherton, 1982] are available for approximating the DF of a given nonlinearity either by numerical integration or by evaluating the DF of an approximating nonlinear characteristic defined, for example, by a quantized characteristic, linear segmented characteristic, or Fourier series. Table 100.3 gives a list of DFs for some commonly used approximations of nonlinear elements. Several of the results are in terms of the DF for an ideal saturation characteristic of unit slope, $N_s(\delta/a)$, defined in Eq. (100.119).
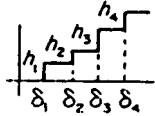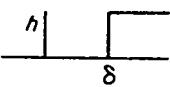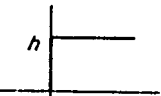
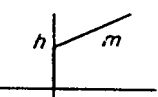**TABLE 100.3**    DFs of Single-Valued Nonlinearities

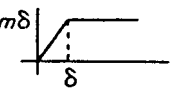| | | |
|---|---|---|
| **General quantizer**  | $a < \delta_1$ <br> $\delta_{M+1} > a > \delta_M$ | $N_p = 0$ <br><br> $N_p = \left(4/a^2\pi\right)\sum\limits_{m=1}^{M} h_m\left(a^2 - \delta_m^2\right)^{1/2}$ |
| **Uniform quantizer** <br> $h_1 = h_2 = \cdots h$ <br> $\delta_m = (2m-1)\delta/2$ | $a < \delta$ <br><br> $(2M+1)\delta > a > (2M-1)\delta$ <br><br> $n = (2m-1)/2$ | $N_p = 0$ <br><br> $N_p = \left(4h/a^2\pi\right)\sum\limits_{m=1}^{M}\left(a^2 - n^2\delta^2\right)^{1/2}$ |
| **Relay with dead zone**  | $a < \delta$ <br> $a > \delta$ | $N_p = 0$ <br> $N_p = 4h(a^2 - \delta^2)^{1/2}/a^2\pi$ |
| **Ideal relay**  | | $N_p = 4h/a\pi$ |
| **Preload**  | | $N_p = (4h/a\pi) + m$ |
| **General piecewise linear**  | $a < \delta_1$ <br> $\delta_{M+1} > \alpha > \delta_M$ | $N_p = (4h/a\pi) + m_1$ <br> $N_p = (4h/a\pi) + m_{M+1}$ <br><br> $+\sum\limits_{i=1}^{M}\left(m_j - m_{j+1}\right)N_s\left(\delta_j/a\right)$ |
| **Ideal saturation**  | | $N_p = mN_s(\delta/a)$ |
| **Dead zone**  | | $N_p = m[1 - N_s(\delta/a)]$ |

**TABLE 100.3 (continued)** DFs of Single-Valued Nonlinearities

**Gain changing nonlinearity**



$$N_p = (m_1 - m_2)N_s(\delta/a) + m_2$$

**Saturation with dead zone**



$$N_p = m[N_s(\delta_2/a) - N_s(\delta_1/a)]$$



$$N_p = -m_1 N_s(\delta_1/a) + (m_1 - m_2)N_s(\delta_2/a) + m_2$$



$a < \delta$     $N_p = 0$

$a > \delta$     $N_p = 4h(a^2 - \delta^2)^{1/2}/a_2\pi + m - mN_s(\delta/a)$



$a < \delta$     $N_p = m_1$

$a > \delta$     $N_p = (m_1 - m_2)N_s(\delta/a) + m_2 + 4h(a^2 - \delta^2)^{1/2}/a^2\pi$



$a < \delta$     $N_p = 4h/a\pi$

$a > \delta$     $N_p = 4h/[a - (a^2 - \delta^2)^{1/2}]/a^2\pi$

**Limited field of view**



$$N_p = (m_1 + m_2)N_s(\delta/a) - m_2 N_s[(m_1 + m_2)\delta/m_2 a]$$



$a < \delta$     $N_p = m_1$

$a > \delta$     $N_p = m_1 N_s(\delta/a) - 4m_1\delta(a^2 - \delta^2)^{1/2}/a^2\pi$

**$y = x^m$**       $m > -2$ $\Gamma$ is the gamma function

$$N_p = \frac{\Gamma(m+1)a^{m-1}}{2^{m-1}\Gamma[(3+m)/2]\Gamma[(1+m)/2]}$$

$$= \frac{2}{\sqrt{\pi}}\frac{\Gamma[(m+2)\,2]a^{m-1}}{\Gamma[(m+3)/2]}$$

**FIGURE 100.46** Nyquist plot showing solution for a limit cycle.

## Limit Cycles and Stability

To investigate the possibility of limit cycles in the autonomous closed loop system of Fig. 100.42, the input to the nonlinearity $n(x)$ is assumed to be a sinusoid so that it can be replaced by the amplitude-dependent DF gain $N(a)$. The open loop gain to a sinusoid is thus $N(a)G(j\omega)$ and, therefore, a limit cycle exists if

$$N(a)G(j\omega) = -1 \qquad (100.121)$$

where $G(j\omega) = G_c(j\omega)G_1(j\omega)$. As in general, $G(j\omega)$ is a complex function of $\omega$ and $N(a)$ is a complex function of $a$, solving Eq. (100.121) will yield both the frequency $\omega$ and amplitude $a$ of a possible limit cycle.

A common procedure to examine solutions of Eq. (100.120) is to use a Nyquist diagram, where the $G(j\omega)$ and $C(a) = -1/N(a)$ loci are plotted as in Fig. 100.46, where they are shown intersecting for $a = a_0$ and $\omega = \omega_0$. The DF method indicates therefore that the system has a limit cycle with the input sinusoid to the nonlinearity, $x$, equal to $a_0 \sin (\omega_0 t + \phi)$, where $\phi$ depends on the initial conditions. When the $G(j\omega)$ and $C(a)$ loci do not intersect, the DF method predicts that no limit cycle will exist if the Nyquist stability criterion is satisfied for $G(j\omega)$ with respect to any point on the $C(a)$ locus. Obviously, if the nonlinearity has unit gain for small inputs, the point $(-1, j0)$ will lie on $C(a)$ and may be used as the critical point, analogous to a linear system.
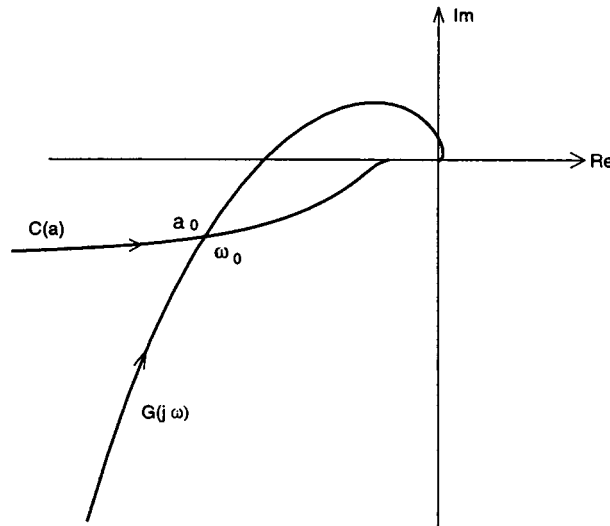
For a stable case, it is possible to use the gain and phase margin to judge the relative stability of the system. However, a gain and phase margin can be found for every amplitude $a$ on the $C(a)$ locus, so it is usually appropriate to use the minimum values of the quantities [Atherton, 1982]. When the nonlinear block includes dynamics so that its response is both amplitude and frequency dependent, that is $N(a, \omega)$, then a limit cycle will exist if

$$G(j\omega) = -1/N(a, \omega) = C(a, \omega). \qquad (100.122)$$

To check for possible solutions of this equation, a family of $C(a, \omega)$ loci, usually as functions of $a$ for fixed values of $\omega$, is drawn on the Nyquist diagram.

An additional point of interest is whether when a solution to Eq. (100.120) exists the predicted limit cycle is stable. When there is only one intersection point, the stability of the limit cycle can be found using the Loeb criterion which states that if the Nyquist stability criterion indicates instability (stability) for the point on $C(a)$ with $a < a_0$ and stability (instability) for the point on $C(a)$ with $a > a_0$ the limit cycle is stable (unstable).

When multiple solutions exist, the situation is more complicated and the criterion above is a necessary but not sufficient result for the stability of the limit cycle [Choudhury and Atherton, 1974].

Normally in these cases, the stability of the limit cycle can be ascertained by examining the roots of the characteristic equation

$$1 + N_{i\gamma}(a)G(s) = 0 \tag{100.123}$$

where $N_{i\gamma}(a)$ is known as the incremental describing function (IDF). $N_{i\gamma}(a)$ for a single valued nonlinearity can be evaluated from

$$N_{i\gamma}(a) = \int_{-a}^{a} n'(x)p(x)dx \tag{100.124}$$

where $n'(x)$ and $p(x)$ are as previously defined. $N_{i\gamma}(a)$ is related to $N(a)$ by the equation

$$N_{i\gamma}(a) = N(a) + (a/2)\,dN(a)/da. \tag{100.125}$$

Thus, for example, for an ideal relay, making $\delta = \Delta = 0$ in Eq. (100.120) gives $N(a) = 4h/a\pi$, also found directly from Eq. (100.116), and, substituting this value in Eq. (100.125) yields $N_{i\gamma}(a) = 2h/a\pi$. Some examples of feedback system analysis using the DF follow.

### Autotuning in Process Control

In 1943 Ziegler and Nichols [1943] suggested a technique for tuning the parameters of a PID controller. Their method was based on testing the plant in a closed loop with the PID controller in the proportional mode. The proportional gain was increased until the loop started to oscillate and then the value of gain and the oscillation frequency were measured. Formulae were given for setting the controller parameters based on the gain named the critical gain, $K_c$, and the frequency called the critical frequency, $\omega_c$.

Assuming that the plant has a linear transfer function $G_1(s)$, then $K_c$ is its gain margin and $\omega_c$ the frequency at which its phase shift is 180°. Performing this test in practice may prove difficult. If the plant has a linear transfer function and the gain is adjusted too quickly, a large amplitude oscillation may start to build up. In 1984 Astrom and Hagglund [1984] suggested replacing the proportional control by a relay element to control the amplitude of the oscillation. Consider therefore the feedback loop of Fig. 100.42 with $n(x)$ an ideal relay, $G_c(s) = 1$, and the plant with a transfer function $G_1(s) = 10/(s + 1)^3$. The $C(a)$ locus, $-1/N(a) = -a\pi/4h$, and the Nyquist locus $G(j\omega)$ in Fig. 100.47 intersect. The values of $a$ and $\omega$ at the intersection can be calculated from

$$-a\pi/4h = \frac{10}{(1 + j\omega)^3} \tag{100.126}$$

which can be written

$$\mathrm{Arg}\left(\frac{10}{(1 + j\omega)^3}\right) = 180°, \text{ and} \tag{100.127}$$

$$\frac{a\pi}{4h} = \frac{10}{(1 + \omega^2)^{3/2}}. \tag{100.128}$$

**FIGURE 100.47** Nyquist plot $10/(s + 1)^3$ and $C(a)$ loci for $\Delta = 0$ and $4h/\pi$.

The solution for $\omega_c$ from Eq. (100.127) is $\tan^{-1} \omega_c = 60°$, giving $\omega_c = \sqrt{3}$. Because the DF solution is approximate, the actual measured frequency of oscillation will differ from this value by an amount which will be smaller the closer the oscillation is to a sinusoid. The exact frequency of oscillation in this case will be 1.708 rads/sec in error by a relatively small amount. For a square wave input to the plant at this frequency, the plant output signal will be distorted by a small percentage. The distortion, $d$, is defined by

$$d = \left[ \frac{\text{M.S. value of signal} \quad - \quad \text{M.S. value of fundamental harmonic}}{\text{M.S. value of fundamental harmonic}} \right]^{1/2} \quad (100.129)$$

Solving Eq. (100.128) gives the amplitude of oscillation $a$ as $5h/\pi$. The gain through the relay is $N(a)$ equal to the critical gain $K_c$. In the practical situation where $a$ is measured, $K_c$ equal to $4h/a\pi$, should be close to the known value of 0.8 for this transfer function.

If the relay has an hysteresis of $\Delta$, then with $\delta = 0$ in Eq. (100.120) gives

$$N(a) = \frac{4h\left(a^2 - \Delta^2\right)^{1/2}}{a^2\pi} - j\frac{4h\Delta}{a^2\pi}$$

from which

$$C(a) = \frac{-1}{N(a)} = \frac{-\pi}{4h}\left[\left(a^2 - \Delta^2\right)^{1/2} + j\Delta\right].$$

Thus on the Nyquist plot, $C(a)$ is a line parallel to the real axis at a distance $\pi\Delta/4h$ below it, as shown in Fig. 100.47 for $\Delta = 1$ and $h = \pi/4$ giving $C(a) = -(a^2 - 1)^{1/2} - j$. If the same transfer function is used for the plant, then the limit cycle solution is given by

$$-\left(a^2 - 1\right)^{1/2} - j = \frac{10}{\left(1 + j\omega\right)^3} \quad (100.130)$$

**FIGURE 100.48**   $N(a)$ for ideal relay with dead zone.



**FIGURE 100.49**   Two limit cycles: $a_1$, unstable; $a_2$, stable.

where $\omega = 1.266$, which compares with an exact solution value of 1.254, and $a = 1.91$. For the oscillation with the ideal relay, Eq. (100.123) with $N_{i\gamma}(a) = 2h/a\pi$ shows that the limit cycle is stable. This agrees with the perturbation approach which also shows that the limit cycle is stable when the relay has hysteresis.

### Feedback Loop with a Relay with Dead Zone

For this example the feedback loop of Fig. 100.42 is considered with $n(x)$ a relay with dead zone and $G(s) = 2/s(s + 1)^2$. From Equation 19.22 with $\Delta = 0$, the DF for this relay, given by

$$N(a) = 4h\left(a^2 - \delta^2\right)^{1/2} \Big/ a^2\pi \text{ for } a > \delta. \tag{100.131}$$

is real because the nonlinearity is single valued. A graph of $N(a)$ against $a$ is in Fig. 100.48, and shows that $N(a)$ starts at zero, when $a = \delta$, increases to a maximum, with a value of $2h/\pi\delta$ at $a = \delta\sqrt{2}$, and then decreases toward zero for larger inputs. The $C(a)$ locus, shown in Fig. 100.49, lies on the negative real axis starting at $-\infty$ and returning there after reaching a maximum value of $-\pi\delta/2h$. The given transfer function $G(j\omega)$ crosses the negative real axis, as shown in Fig. 100.49, at a frequency of $\tan^{-1}\omega = 45°$, that, is $\omega = 1$ rad/sec and, therefore, cuts the $C(a)$ locus twice. The two possible limit cycle amplitudes at this frequency can be found by solving

$$\frac{a^2\pi}{4h\left(a^2 - \delta^2\right)^{1/2}} = 1$$

**FIGURE 100.50**  Circle criterion and stability.

which gives $a = 1.04$ and 3.86 for $\delta = 1$ and $h = \pi$. Using the perturbation method or the IDF criterion, the smallest amplitude limit cycle is unstable and the larger one is stable. If a condition similar to the lower amplitude limit cycle is excited in the system, an oscillation will build up and stabilize at the higher amplitude limit cycle.
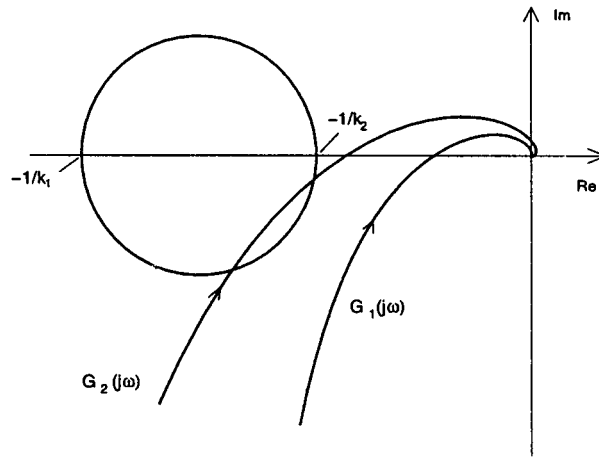
Other techniques show that the exact frequencies of the limit cycles for the smaller and larger amplitudes are 0.709 and 0.989, respectively. Although the transfer function is a good low pass filter, the frequency of the smallest amplitude limit cycle is not predicted accurately because the output from the relay, a wave form with narrow pulses, is highly distorted.

If the transfer function of $G(s)$ is $K/s(s + 1)^2$, then no limit cycle will exist in the feedback loop, and it will be stable if

$$\left.\frac{K}{\omega\left(1 + \omega^2\right)}\right|_{\omega=1} < \frac{\pi d}{2h} \, ,$$

that is, $K < \pi\delta/h$. If $\delta = 1$ and $h = \pi$, $K < 1$ which may be compared with the exact result for stability of $K < 0.96$.

## Stability and Accuracy

Because the DF method is an approximate procedure, it is desirable to judge its accuracy. Predicting that a system will be stable, when in practice it is not, may have unfortunate consequences. Many attempts have been made to solve this problem, but those obtained are difficult to apply or produce too conservative results [Mess and Bergen, 1975].

The problem is illustrated by the system of Fig. 100.42 with a symmetrical odd single-valued nonlinearity confined to a sector between lines of slope $k_1$ and $k_2$, that is, $k_1 x < n(x) < k_2 x$ for $x > 0$. For absolute stability, the circle criterion requires satisfying the Nyquist criterion for the locus $G(j\omega)$ for all points within a circle having its diameter on the negative real axis of the Nyquist diagram between the points $(-1/k_1, 0)$ and $(-1/k_2, 0)$, as shown in Fig. 100.50. On the other hand, because the DF for this nonlinearity lies within the diameter of the circle, the DF method requires satisfying the Nyquist criterion for $G(j\omega)$ for all points on the circle diameter, if the autonomous system is to be stable.

Therefore, for a limit cycle in the system of Fig. 100.42, errors in the DF method relate to its inability to predict a phase shift, which the fundamental harmonic may experience in passing through the nonlinearity, rather than an incorrect magnitude of the gain. When the input to a single-valued nonlinearity is a sinusoid together with some of its harmonics, the fundamental output is not necessarily in phase with the fundamental

input, that is, the fundamental gain has a phase shift. The actual phase shift varies with the harmonic content of the input signal in a complex manner, because the phase shift depends on the amplitudes and phases of the individual input components.

From an engineering viewpoint one can judge the accuracy of DF results by estimating the distortion, *d,* in the input to the nonlinearity. This is straightforward when a limit-cycle solution is given by the DF method; the loop may be considered opened at the nonlinearity input, the sinusoidal signal corresponding to the DF solution can be applied to the nonlinearity, and the harmonic content of the signal fed back to the nonlinearity input can be calculated. Experience indicates that the percentage accuracy of the DF method in predicting the fundamental amplitude and frequency of the limit cycle is less than the percentage distortion in the fedback signal. As mentioned previously, the DF method may incorrectly predict stability. To investigate this problem, the procedure above can be used again, by taking, as the nonlinearity input, a sinusoid with amplitude and frequency corresponding to values of those parameters where the phase margin is small. If the calculated fedback distortion is high, say greater than 2% per degree of phase margin, the DF result should not be relied on.

The limit-cycle amplitude predicted by the DF is an approximation to the fundamental harmonic. The accuracy of this prediction cannot be assessed by using the peak value of the limit cycle to estimate an equivalent sinusoid. It is possible to estimate the limit cycle more accurately by balancing more harmonics, as mentioned earlier. Although this is difficult algebraically other than with loops whose nonlinearity is mathematically simply described, for example a cubic, software is available for this purpose [McNamara and Atherton, 1987]. The procedure involves solving sets of nonlinear algebraic equations but good starting guesses can usually be obtained for the magnitudes and phases of the other harmonic components from the wave form fedback to the nonlinearity, assuming its input is the DF solution.

## Compensator Design

Although the design specifications for a control system are often in terms of step-response behavior, frequency domain design methods rely on the premise that the correlation between the frequency and a step response yields a less oscillatory step response if the gain and phase margins are increased. Therefore the design of a suitable linear compensator for the system of Fig. 100.42 using the DF method, is usually done by selecting for example a lead network to provide adequate gain and phase margins for all amplitudes. This approach may be used in example 2 of the previous section where a phase lead network could be added to stabilize the system, say for a gain of 1.5, for which it is unstable without compensation. Other approaches are the use of additional feedback signals or modification of the nonlinearity $n(x)$ directly or indirectly [Atherton, 1982; Gelb and van der Velde, 1968].

When the plant is nonlinear, its frequency response also depends on the input sinusoidal amplitude represented as $G(j\omega, a)$. In recent years several approaches [Nanka-Bruce and Atherton, 1990; Taylor and Strobel, 1984] use the DF method to design a nonlinear compensator for the plant, with the objective of closed-loop performance independent of the input amplitude.

## Closed-Loop Frequency Response

When the closed-loop system of Fig. 100.42 has a sinusoidal input $r(t) = R\sin(\omega t + \theta)$, it is possible to evaluate the closed-loop frequency response using the DF. If the feedback loop has no limit cycle when $r(t) = 0$ and, in addition, the sinusoidal input $r(t)$ does not induce a limit cycle, then, provided that $G_c(s)G_1(s)$ gives good filtering, $x(t)$, the nonlinearity input, almost equals the sinusoid $a\sin\omega t$. Balancing the components of frequency $\omega$ around the loop,

$$g_c R \sin(\omega t + \theta - \phi_c) - ag_1 g_c M(a)$$

$$\sin[\omega t + \phi_1 + \phi_c + \psi(a)] = a \sin \omega t \qquad (100.132)$$

where $G_c(j\omega) = g_c e^{j\phi_c}$ and $G_1(j\omega) = g_1 e^{j\phi_1}$. In principle Eq. (100.132), which can be written as two nonlinear algebraic equations, can be solved for the two unknowns $a$ and $\theta$ and the fundamental output signal can then be found from

$$c(t) = aM(a)g_1 \sin\left[\omega t + \psi(a) + \phi_1\right] \tag{100.133}$$

to obtain the closed-loop frequency for $R$ and $\omega$.

Various graphical procedures have been proposed for solving the two nonlinear algebraic equations resulting from Eq. (100.132) [Levinson, 1953; Singh, 1965; West and Douce, 1954]. If the system is lightly damped, the nonlinear equations may have more than one solution, indicating that the frequency response of the system has a jump resonance. This phenomenon of a nonlinear system has been studied by many authors, both theoretically and practically [Lamba and Kavanagh, 1971; West et al., 1954].

## The Phase Plane Method

The phase plane method was the first method used by control engineers for studying the effects of nonlinearity in feedback systems. The technique which can only be used for systems with second order models was examined and further developed for control engineering purposes for several major reasons,

1. The phase plane approach has been used for several studies of second order nonlinear differential equations arising in fields such as planetary motion, nonlinear mechanics and oscillations in vacuum tube circuits.
2. Many of the control systems of interest, such as servomechanisms, could be approximated by second order nonlinear differential equations.
3. The phase plane was particularly appropriate for dealing with nonlinearities with linear segmented characteristics which were good approximations for the nonlinear phenomena encountered in control systems.

The next section considers the basic aspects of the phase plane approach but later concentration is focused on control engineering applications where the nonlinear effects are approximated by linear segmented nonlinearities.

### Background

Early analytical work [Andronov et al., 1966], on second order models assumed the equations

$$\begin{aligned}
\dot{x}_1 &= P(x_1, x_2) \\
\dot{x}_2 &= Q(x_1, x_2)
\end{aligned} \tag{100.134}$$

for two first-order nonlinear differential equations. Equilibrium, or singular points, occur when

$$\dot{x}_1 = \dot{x}_2 = 0$$

and the slope of any solution curve, or trajectory, in the $x_1 - x_2$ state plane is

$$\frac{dx_2}{dx_1} = \frac{\dot{x}_2}{\dot{x}_1} = \frac{Q(x_1, x_2)}{P(x_1, x_2)} \tag{100.135}$$

A second order nonlinear differential equation representing a control system can be written

$$\ddot{x} + f(x, \dot{x}) = 0 \qquad (100.136)$$

If this is rearranged as two first-order equations, choosing the phase variables as the state variables, that is $x_1 = x$, $x_2 = \dot{x}$, then Eq. (100.136) can be written as

$$\dot{x}_1 = \dot{x}_2 \qquad \dot{x}_2 = -f(x_1, x_2) \qquad (100.137)$$

which is a special case of Eq. (100.135). A variety of procedures has been proposed for sketching state [phase] plane trajectories for Eqs. (100.135) and (100.137). A complete plot showing trajectory motions throughout the entire state (phase) plane is known as a state (phase) portrait. Knowledge of these methods, despite the improvements in computation since they were originally proposed, can be particularly helpful for obtaining an appreciation of the system behavior. When simulation studies are undertaken, phase plane graphs are easily obtained and they are often more helpful for understanding the system behavior than displays of the variables $x_1$ and $x_2$ against time.

Many investigations using the phase plane technique were concerned with the possibility of limit cycles in the nonlinear differential equations When a limit cycle exists, this results in a closed trajectory in the phase plane. Typical of such investigations was the work of Van der Pol, who considered the equation

$$\ddot{x} - \mu(1 - x^2)\dot{x} + x = 0 \qquad (100.138)$$

where $\mu$ is a positive constant. The phase plane form of this equation can be written as

$$\begin{aligned} \dot{x}_1 &= x_2 \\ \dot{x}_2 &= -f(x_1, x_2) = \mu(1 - x_1^2)x_2 - x_1 \end{aligned} \qquad (100.139)$$

The slope of a trajectory in the phase plane is

$$\frac{dx_2}{dx_1} = \frac{\dot{x}_2}{\dot{x}_1} = \frac{\mu(1 - x_1^2)x_2 - x_1}{x_2} \qquad (100.140)$$

and this is only singular (that is, at an equilibrium point), when the right hand side of Eq. (100.140) is 0/0, that is $x_1 = x_2 = 0$.

The form of this singular point which is obtained from linearization of the equation at the origin depends upon $\mu$, being an unstable focus for $\mu < 2$ and an unstable node for $\mu > 2$. All phase plane trajectories have a slope of $r$ when they intersect the curve

$$rx_2 = \mu(1 - x_1^2)x_2 - x_1 \qquad (100.141)$$

One way of sketching phase plane behavior is to draw a set of curves given for various values of $r$ by Eq. (100.141) and marking the trajectory slope r on the curves. This procedure is known as the method of isoclines and has been used to obtain the limit cycles shown in Fig. 100.51 for the Van der Pol equation with $\mu = 0.2$ and 4.

## Piecewise Linear Characteristics

When the nonlinear elements occurring in a second order model can be approximated by linear segmented characteristics then the phase plane approach is usually easy to use because the nonlinearities divide the phase

**FIGURE 100.51**    Phase portraits of the Van der Pol equation for different values, of $\mu$.

plane into various regions within which the motion may be described by different linear second-order equations [Atherton, 1982]. The procedure is illustrated by the simple relay system in Fig. 100.52.

The block diagram represents an "ideal" relay position control system with velocity feedback. The plant is a double integrator, ignoring viscous (linear) friction, hysteresis in the relay, or backlash in the gearing. If the system output is denoted by $x_1$ and its derivative by $x_2$, then the relay switches when $-x_1 - x_2 = \pm 1$; the equations of the dotted lines are marked switching lines on Fig. 100.53.

Because the relay output provides constant values of $\pm 2$ and 0 to the double integrator plant, if we denote the constant value by $h$, then the state equations for the motion are

**FIGURE 100.52**  Relay system.



**FIGURE 100.53**  Phase plane for relay system.

$$\dot{x}_1 = x_2$$
$$\dot{x}_2 = h \qquad (100.142)$$

which can be solved to give the phase plane equation

$$x_2^2 - x_{20}^2 = 2h(x_1 - x_{10}) \qquad (100.143)$$

which is a parabola for $h$ finite and the straight line $x_2 = x_{20}$ for $h = 0$, where $x_{20}$ and $x_{10}$ are the initial values of $x_2$ and $x_1$. Similarly, more complex equations can be derived for other second-order transfer functions. Using Eq. (100.143) with the appropriate values of $h$ for the three regions in the phase plane, the step response for an input of 4.6 units can be obtained as shown in Fig. 100.53.

In the step response, when the trajectory meets the switching line $x_1 + x_2 = -1$ for the second time, trajectory motions at both sides of the line are directed towards it, resulting in a sliding motion down the switching line. Completing the phase portrait by drawing responses from other initial conditions shows that the autonomous system is stable and also that all responses will finally slide down a switching line to equilibrium at $x_1 = \pm 1$.

An advantage of the phase plane method is that it can be used for systems with more than one nonlinearity and for those situations where parameters change as functions of the phase variables. For example, Fig. 100.54 shows the block diagram of an approximate model of a servomechanism with nonlinear effects due to torque saturation and Coulomb friction.

**FIGURE 100.54**  Block diagram of servomechanism.

The differential equation of motion in phase variable form is

$$\dot{x}_2 \;=\; f_s\!\left(-x_1\right) - \left(1/2\right) \operatorname{sgn} x_2 \qquad\qquad (100.144)$$

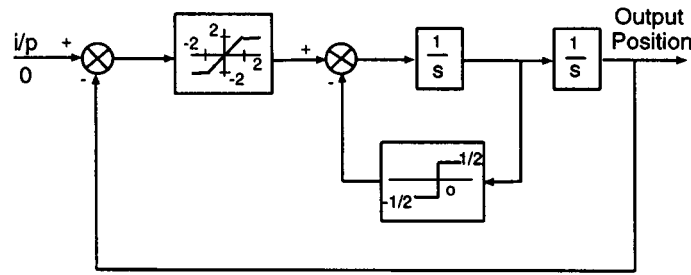where $f_s$ denotes the saturation nonlinearity and sgn the signum function, which is +1 for $x_2 > 0$ and −1 for $x_2 < 0$. There are six linear differential equations describing the motion in different regions of the phase plane. For $x_2$ positive, Eq. (100.144) can be written

$$\dot{x}_1 \;+\; f_s\!\left(x_1\right) + 1/2 \;=\; 0$$

so that for

   (a)  $x_2$+ve, $x_1 < -2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 = 3/2$, a parabola in the phase plane.
   (b)  $x_2$+ve$|x_1| < 2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 + x_1 + 1/2 = 0$.
   (c)  $x_2$+ve, $x_1 > 2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 = -5/2$, a parabola in the phase plane. Similarly for $x_2$ negative,
   (d)  $x_2$−ve, $x_1 − 2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 = -5/2$, a parabola in the phase plane.
   (e)  $x_2$−ve, $|x_2| < 2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 + x_1 - 1/2 = 0$, a circle in the phase plane.
   (f)  $x_2$−ve, $x_1 > 2$, we have $\dot{x}_1 = x_2$, $\dot{x}_2 = -3/2$, a parabola in the phase plane.

Because all the phase plane trajectories are described by simple mathematical expressions, it is straightforward to calculate specific phase plane trajectories.

## Discussion

The phase plane approach is useful for understanding the effects of nonlinearity in second order systems, particularly if it may be approximated by a linear segmented characteristic. Solutions for the trajectories with other nonlinear characteristics may not be possible analytically so that approximate sketching techniques were used in early work on nonlinear control. These approaches are described in many books, for example, [Blaquiere, 1966; Cosgriff, 1958; Cunningham, 1958; Gibson, 1963; Graham and McRuer, 1961; Hayashi, 1964, Thaler and Pastel, 1962; West, 1960]. Although the trajectories are now easily obtained with modern simulation techniques, knowledge of the topological aspects of the phase plane are still useful for interpreting the responses in different regions of the phase plane and appreciating the system behavior.

## Related Topics

5.2 Limiters  •  12.1 Introduction  •  12.3 Lyapunov Stability Theory

## References

A.A. Andronov, A.A. Vitt, and S.E. Khaikin, *Theory of Oscillators,* Reading, Mass.: Addison-Wesley, 1966. (First edition published in Russia in 1937.)

K.J. Astrom, and T. Haggland, *Automatic tuning of single regulators,* Budapest: Proc IFAC Congress, Vol. 4, 267–272, 1984.

D.P. Atherton, *Nonlinear Control Engineering, Describing Function Analysis and Design,* London: Van Nostrand Reinhold, 1975.

D.P. Atherton, *Non Linear Control Engineering,* Student Ed., New York: Van Nostrand Reinhold, 1982.

A. Blaquiere, *Nonlinear Systems Analysis,* New York: Academic Press, 1966.

S.K. Choudhury, and D.P. Atherton, "Limit cycles in high order nonlinear systems," *Proc. Inst. Electr. Eng.,* 121, 717–724, 1974.

P.A. Cook, "Describing function for a sector nonlinearity," *Proc. Inst. Electr. Eng.,* 120, 143–144, 1973.

R. Cosgriff, *Nonlinear Control Systems,* New York: McGraw-Hill, 1958.

W.J. Cunningham, *Introduction to Nonlinear Analysis,* New York: McGraw-Hill, 1958.

A. Gelb and W.E. van der Velde, *Multiple Input Describing Functions and Nonlinear Systems Design,* New York: McGraw-Hill, 1968.

J.E. Gibson, *Nonlinear Automatic Control,* New York: McGraw-Hill, 1963.

D. Graham and D. McRuer, *Analysis of Nonlinear Control Systems,* New York: John Wiley & Sons, 1961.

C. Hayashi, *Nonlinear Oscillations in Physical Systems,* New York, McGraw-Hill, 1964.

S.S. Lamba and R.J. Kavanagh, "The phenomenon of isolated jump resonance and its application," *Proc. Inst. Electr. Eng.,* 118, 1047–1050, 1971.

E. Levinson, "Some saturation phenomena in servomechanims with emphasis on the techometer stabilised system," *Trans, Am. Inst. Electr. Eng.,* Part 2, 72, 1–9, 1953.

O.P. McNamara, and D.P. Atherton, "Limit cycle prediction in free structured nonlinear systems," *IFAC Congress,* Munich, 8, 23–28, July 1987.

A.I. Mees and A.R. Bergen, "Describing function revisited," *IEEE Trans. Autom. Control,* 20, 473–478, 1975.

O. Nanka-Bruce and D.P. Atherton, "Design of nonlinear controllers for nonlinear plants," *IFAC Congress,* Tallinn, 6, 75–80, 1990.

T.P. Singh, "Graphical method for finding the closed loop frequency response of nonlinear feedback control systems," *Proc. Inst. Electr. Eng.,* 112, 2167–2170, 1965.

J.H. Taylor and K.L. Strobel, "Applications of a nonlinear controller design approach based on the quasilinear system models," *Prof ACC,* San Diego, 817–824, 1984.

G.J. Thaler and M.P. Pastel, *Analysis and Design of Nonlinear Feedback Control Systems,* New York: McGraw-Hill, 1962.

J.C. West, *Analytical Techniques of Nonlinear Control Systems,* London: E.U.P., 1960.

J.C. West and J.L. Douce, "The frequency response of a certain class of nonlinear feedback systems," *Br. J. Appl. Phys.,* 5, 201–210, 1954.

J.C. West, B.W. Jayawant, and D.P. Rea, "Transition characteristics of the jump phenomenon in nonlinear resonant circuits," *Proc. Inst. Electr. Eng.,* 114, 381–392, 1967.

J.G. Ziegler and N.B. Nichols, "Optimal setting for automatic controllers," *Trans. ASME,* 65, 433–444, 1943.

## Further Information

Many control engineering text books contain material on nonlinear systems where the describing function is discussed. The coverage, however, is usually restricted to the basic sinusoidal DF for determining limit cycles in feedback systems. The basic DF method, which is one of quasilinearisation, can be extended to cover other signals, such as random signals, and also to cover multiple input signals to nonlinearities and feedback system analysis. The two books with the most comprehensive coverage of this are Gelb and Van der Velde [1968] and Atherton [1975]. More specialized books on nonlinear feedback systems usually cover the phase plane method and the DF, together with other topics such as absolute stability, exact linearization, etc.

## 100.8 Optimal Control and Estimation

*John S. Bay and William T. Baumann*

Consider the closed-loop feedback control of linear time-invariant, multi-input/multi-output (MIMO) state-space systems of the form:

$$\dot{x}(t) = Ax(t) + Bu(t)$$
$$y(t) = Cx(t) + Du(t)$$

(100.145)

In this form, the vector $x \in \Re^n$ represents the internal state, $u \in \Re^p$ represents the input, and $y \in \Re^q$ represents the measured outputs. It is well known that if the system in Eq. (100.145) is *stabilizable* (i.e., its unstable part is **controllable**), then it can be asymptotically stabilized with static state feedback. If it is **detectable** (i.e., its unstable part is **observable**), then a state estimator can be found, whose state variables asymptotically approach the true state variables in Eq. (100.145). However, merely determining the state feedback gain or observer gain leaves considerable design freedom for satisfying criteria other than stabilization and asymptotic observation. In this chapter section, we will provide results of some basic techniques of optimal control and estimation, which provide a mechanism to find *optimal* feedback and **observer (estimator)** gains according to selected optimality criteria.

### Linear Quadratic Regulators

The linear quadratic regulator (LQR) problem is to find an *optimal* control input $u^*(t)$ that minimizes the performance criterion

$$J(x,u) = \tfrac{1}{2} x^{\mathrm{T}}(t_f) sx(t_f) + \tfrac{1}{2} \int_{t_o}^{t_f} \left[ x(t) Qx(t) + u^{\mathrm{T}}(t) Ru(t) \right] dt$$

(100.146)

where $S$ and $Q$ are symmetric, **positive-semidefinite** weighting matrices; and $R$ is a symmetric, *positive-definite* weighting matrix. In this criterion, the term $\tfrac{1}{2} x^{\mathrm{T}}(t_f) Sx(t_f)$ represents a penalty for the state at the final time $t_f$ being different from zero. The term inside the integral, $x^{\mathrm{T}}(t) Qx(t)$, represents a penalty on the transient response of the state vector. The term $u^{\mathrm{T}}(t) Ru(t)$ represents a penalty on the size of the control input $u(t)$. We allow $S$ and $Q$ to be positive-semidefinite because we can generally tolerate unbounded state variables, provided they are not observed at the output. However, by forcing $R$ to be positive-definite, we can guarantee that the process of minimizing Eq. (100.146) gives a bounded input. Minimization of this control energy is one of the primary reasons for using optimal control.

The optimal control $u^*(t)$ can be found via a number of techniques, including dynamic programming [Bay, 1999] and variational techniques [Kirk, 1970]. The result of any of these methods is that the optimal control $u^*(t)$ is a linear function of the state vector (linear state feedback) of the form:

$$u^*(t) = -R^{-1} B^{\mathrm{T}} P(t) x(t)$$

(100.147)

where the $n \times n$ matrix function $P(t)$ satisfies the following equation:

$$\dot{P} = PBR^{-1} B^{\mathrm{T}} P - Q - PA - A^{\mathrm{T}} P$$

(100.148)

Equation (100.148) is known as the differential matrix Riccati equation, and it is solved in backward time, with the end-time condition $P(t_f) = S$.

It may be noted that a reasonable optimization criterion may have no finite final time $t_f$. Instead, it may be desired that the controller be continually active, implying $t_f \to \infty$ and eliminating the possibility of the final state term in Eq. (100.146). In this case, the optimization criterion is more properly written as

$$J(x,u) = \tfrac{1}{2} \int_{t_o}^{\infty} \left[ x^{\mathrm{T}}(t) Q x(t) + u^{\mathrm{T}}(t) R u(t) \right] dt \qquad (100.149)$$

Fortunately, this criterion simplifies the optimal control solution to the steady-state solution of the finite-time problem. That is, for the criterion of Eq. (100.149), the optimal control is

$$u^*(t) = -R^{-1} B^{\mathrm{T}} P x(t) \qquad (100.150)$$

where in this case $P$ is the matrix solution to the following algebraic Riccati equation:

$$0 = P B R^{-1} B^{\mathrm{T}} P - Q - PA - A^{\mathrm{T}} P \qquad (100.151)$$

Such a steady-state optimal solution exists whenever the system $(A, B)$ is stabilizable. Furthermore, this constant $P$ is the unique positive-definite solution of Eq. (100.151) if and only if the pair $(A, T)$ is detectable, where $T$ is defined as the square-root of $Q$, $Q = T^{\mathrm{T}}T$. **Stabilizability** of $(A, B)$ ensures convergence of the cost criterion integral Eq. (100.149), and detectability of $(A, T)$ guarantees that no unstable part of $x(t)$ escapes integration as part of the integrand of Eq. (100.149).

We should point out also that for the corresponding discrete-time system:

$$\begin{aligned} x(k+1) &= Ax(k) + Bu(k) \\ y(k) &= Cx(k) + Du(k) \end{aligned} \qquad (100.152)$$

minimization of the cost criterion:

$$J = \tfrac{1}{2} \sum_{k=k_0}^{\infty} \left[ x^{\mathrm{T}}(k) Q x(k) + u^{\mathrm{T}}(k) R u(k) \right] \qquad (100.153)$$

over all inputs $u(k)$ results in the optimal control

$$u^*(k) = -\left[ R + B^{\mathrm{T}} SB \right]^{-1} B^{\mathrm{T}} SAx(k) \qquad (100.154)$$

where $S$ is the solution to the corresponding discrete-time algebraic Riccati equation:

$$S = A^{\mathrm{T}} SA - A^{\mathrm{T}} SB \left[ R + B^{\mathrm{T}} SB \right]^{-1} B^{\mathrm{T}} SA + Q \qquad (100.155)$$

Note that in both the continuous- and the discrete-time infinite-horizon cases, the optimal control is actually static state feedback of the form $u = Kx$.

## Optimal Estimation: The Kalman Filter

It was noted above that the optimal controller for the linear quadratic cost criterion takes the form of full-state feedback. However, it is often the case that the full state is not physically available for feedback. Rather, it is usually the *output* that is measurable, so that we prefer a technique that uses $y(t)$ (and possibly $u(t)$) to construct the control signal instead of the state $x(t)$.

The simple solution to this problem is to design an **_observer (or estimator)_**, which produces an *estimated* state vector $\hat{x}(t)$. If this estimate asymptotically approaches the true state $x(t)$, then we can simply combine the observer and the state feedback to produce a feedback control $u(t) = K\hat{x}(t)$. That we can simply substitute the observed value $\hat{x}(t)$ for $x(t)$ in the feedback function is a fortunate property called the separation principle, which ensures that our controller calculations and observer calculations do not interfere with one another.

Just as we have improved static state feedback by introducing the optimal LQR above, we can take the principles of observer design and extend them with some guarantees of optimality. Such an optimal estimator is the Kalman filter.

The Kalman filter is derived assuming the system model:

$$\dot{x}(t) = Ax(t) + Bu(t) + Gv(t)$$
$$y(t) = Cx(t) + w(t)$$

(100.156)

where $v(t)$ and $w(t)$ are two white, Gaussian, zero-mean, mutually uncorrelated noise signals with $E[w(t)v^{\mathrm{T}}(\tau)] = 0$ and

$$E\big[v(t)\big] = 0, \quad E\big[v(t)v^{\mathrm{T}}(\tau)\big] = V\delta(t - \tau)$$

(100.157)

and

$$E\big[w(t)\big] = 0, \quad E\big[w(t)w^{\mathrm{T}}(\tau)\big] = W\delta(t - \tau)$$

(100.158)

where $\delta(t)$ is the Dirac delta. Noise $v(t)$ is called the *plant noise*, and $w(t)$ is the *measurement noise*, often representing sensor noise. (By assuming these signals are first passed through auxiliary filters with specified dynamics, these noises can be made to resemble harmonic or narrow-band disturbances.)

The goal is now to design a system that produces an estimated state $\hat{x}(t)$ for Eq. (100.156) while rejecting the influence of the signals $v(t)$ and $w(t)$. To do this, we need the further assumptions that the system's initial state is guessed to be $x(t_0) = x_0$ and that this guess is uncorrelated with the plant and measurement noise:

$$E\big[x_0 v^{\mathrm{T}}(\tau)\big] = 0 \qquad E\big[x_0 w^{\mathrm{T}}(\tau)\big] = 0$$

(100.159)

The covariance of the initial guess is defined as

$$E\left\{\big[x_0 - E(x_0)\big]\big[x_0 - E(x_0)\big]^{\mathrm{T}}\right\} \triangleq P_0$$

(100.160)

The Kalman filter is the system that performs this estimation, rejecting the influence of the noise. The filter itself is given by the equation:

$$\dot{\hat{x}}(t) = A\hat{x}(t) + Bu(t) + L(t)\big[y(t) - C\hat{x}(t)\big]$$

(100.161)

which can be seen to resemble the standard full-order observer [Bay, 1999]. However, rather than choosing an observer gain $L$ in Eq. (100.161) to simply stabilize the error dynamics of the observer, the *Kalman gain* $L(t)$ is

$$L(t) = P(t)C^\mathrm{T}W^{-1} \qquad (100.162)$$

where $P(t)$ is the solution to the following differential Riccati equation:

$$\dot{P}(t) = AP(t) + P(t)A^\mathrm{T} - P(t)C^\mathrm{T}W^{-1}CP(t) + GVG^\mathrm{T} \qquad (100.163)$$

whose initial condition is $P(t_0) = P_0$.

For the discrete-time system

$$\begin{aligned}
x(k+1) &= Ax(k) + Bu(k) + Gv(k) \\
y(k) &= Cx(k) + w(k)
\end{aligned} \qquad (100.164)$$

with assumptions analogous to Eq. (100.157) through (100.158) for plant noise $v(k)$ and measurement noise $w(t)$, and with $E[(x_0 - E(x_0)(x_0 - E(x_0))^\mathrm{T}] \triangleq S_0$, the Kalman filter is given by the two-stage estimator

$$\bar{x}(k+1) = A\hat{x}(k) + Bu(k) \qquad (100.165)$$

$$\hat{x}(k+1) = \bar{x}(k+1) + L(k+1)\left[ y(k+1) - C\bar{x}(k+1) \right] \qquad (100.166)$$

where the Kalman gain $L(k+1)$ is computed from the equation

$$L(k+1) = \left[ AS(k)A^\mathrm{T} + GVG^\mathrm{T} \right]C^\mathrm{T}\left\{ C\left[ AS(k)A^\mathrm{T} + GVG^\mathrm{T} \right]C^T + W \right\}^{-1} \qquad (100.167)$$

In Eq. (100.167), the term $S(k)$ is determined from the Riccati equation:

$$\begin{aligned}
S(k+1) &= \left[ I - L(k+1)C \right]\left[ AS(k)A^\mathrm{T} + GVG^\mathrm{T} \right]\left[ I - L(k+1)C \right]^\mathrm{T} \\
&\quad + L(k+1)WL^\mathrm{T}(k+1)
\end{aligned} \qquad (100.168)$$

with initial condition $S(k_0) = S_0$. (Note that these equations can be combined and rearranged to produce a number of alternate formulations.)

Equation (100.165) is often referred to as the *time update* equation. It represents the estimate of the state vector that results from knowledge of the system dynamics. Equation (100.166) is sometimes called the *measurement* update equation because it revises the time update with a so-called *innovations* term $L(y - \bar{y})$ that adjusts this time update according to the error between the output expected from the time update, $\bar{y}(k+1)$, and the actual, measured output, $y(k+1)$.

The matrices $P(t)$ in the continuous-time filter, and $S(k)$ in the discrete-time filter are the error covariance matrices for the state estimate. That is,

$$S(k) \triangleq E\left[ e(k)e^\mathrm{T}(k) \right] \quad \text{and} \quad P(t) \triangleq E\left[ e(t)e^\mathrm{T}(t) \right] \qquad (100.169)$$

where $e(k) \triangleq x(k) - \hat{x}(k)$ and $e(t) \triangleq x(t) - \hat{x}(t)$. Thus, the size of these matrices is an indicator of the error variance in various components in the estimates, and it can be seen in Eq. (100.161) and (100.167) that as these covariances decrease, the estimators rely less and less on the innovations term and more on the system dynamics for an accurate estimate. Early in the estimation process, the situation is usually reversed, with the innovations having the larger effect.

## Linear-Quadratic-Gaussian (LQG) Control

It can be shown that the Kalman filter is the *optimal* estimator for the state of system Eq. (100.156) or (100.164) in the sense that it minimizes the squared error due to the noise input terms. Therefore, it becomes a likely candidate for combination with the LQR of the previous section. Together, the combination of LQR and Kalman filter is known as the LQG (linear-quadratic-Gaussian) controller, and is a useful controller in many situations. This controller is optimal in the sense that it minimizes the expected root-mean-square value of the optimization criterion

$$\lim_{T \to \infty} E \left\{ \frac{1}{T} \int_0^T x^{\mathrm{T}}(t) Q x(t) + u^{\mathrm{T}}(t) R u(t) dt \right\}^{1/2} \tag{100.170}$$

when the noise inputs are unit variance white noise. In the frequency domain, this is equivalent to minimizing the $H_2$-norm

$$\|G\|_2 = \left\{ \frac{1}{2\pi} \int_{-\infty}^{\infty} trace\left( G^\star(j\omega) G(j\omega) \right) d\omega \right\}^{1/2} \tag{100.171}$$

of the transfer function $G$ from the white noise inputs $\begin{bmatrix} v \\ w \end{bmatrix}$ to the output $z = \begin{bmatrix} Tx \\ R^{1/2} u \end{bmatrix}$ where $Q = T^{\mathrm{T}} T$, as before.

We should point out that the so-called $H_2$ controller is equivalent to the LQG controller but provides a unified framework for control and observation that explains the striking similarity between the LQ controller equations and the Kalman filter equations (for example, in the Riccati equations of (100.148) and (100.63)). See Zhou and Doyle [1998] for further information.

## $H_\infty$ Control

The standard $H_\infty$ control problem considers a system of the form

$$\dot{x} = Ax + B_1 w + B_2 u$$
$$z = C_1 x + D_{12} u \tag{100.172}$$
$$y = C_2 x + D_{21} w$$

where $w$ is a deterministic disturbance input vector, $y$ is the measured output vector, and $z$ is a vector of variables to be controlled. The objective of $H_\infty$ control is to minimize the $H_\infty$ norm

$$\|G\|_\infty = \sup_\omega \overline{\sigma}\left[ G(j\omega) \right] \tag{100.173}$$

(where $\overline{\sigma}$ denotes the maximum *singular value* of a matrix, and sup denotes "supremum," or least upper bound) of the transfer function $G$ from the disturbance input $w$ to the output $z$. In the time domain, the square of this

norm corresponds to the maximum possible energy magnification between the input and the output, where energy is defined as the integral of the square of a signal; for example, $\int_0^\infty w^T(t)w(t)dt$.

One of the major reasons for the development of $H_\infty$ control is that many performance and robustness criteria for MIMO systems involve the maximum **singular value** of certain system transfer functions. To optimize the performance or robustness of these systems requires the minimization of the maximum singular value of these transfer functions, which is exactly the objective of $H_\infty$ control.

From a disturbance rejection point of view, the $H_\infty$ controller can be used to minimize the root-mean-square value of the controlled variable $z$ due to the worst-case unit-energy disturbance $w$. This is to be contrasted with LQG (or $H_2$) controller, which minimizes the average response to a unit-variance random disturbance.

In practice, it is common to solve the suboptimal $H_\infty$ problem where it is desired to find an output feedback controller such that $\|G\|_\infty < \gamma$, where $\gamma$ is specified by the designer. For large values of $\gamma$, there will always be a solution to the problem. In fact, as $\gamma$ approaches infinity in the equations below, the central $H_\infty$ controller will approach the LQG controller. By decreasing the value of $\gamma$ until just before a solution to the problem ceases to exist, the designer can get as close to the optimal $H_\infty$ controller as desired. To ensure that a solution for some value of $\gamma$ exists, the following standard assumptions on the system are made [Green and Limebeer, 1995]:

1. The pair $(A, B_2)$ is stabilizable and the pair $(A, C_2)$ is detectable
2. $D_{12}^T D_{12} = I$ and $D_{21}D_{21}^T = I$

3. Rank $\begin{bmatrix} A - j\omega I & B_2 \\ C_1 & D_{12} \end{bmatrix} = n + m$    for all real $\omega$

4. Rank $\begin{bmatrix} A - j\omega I & B_1 \\ C_2 & D_{21} \end{bmatrix} = n + q$    for all real $\omega$

where $n$ is the dimension of $x$, $m$ is the dimension of $u$, and $q$ is the dimension of $y$.

Under these assumptions, it can be shown that there exists a stabilizing, measurement feedback solution to the suboptimal $H_\infty$ control problem if and only if the following three conditions are met.

1. The algebraic Riccati equation $X_\infty \tilde{A} + \tilde{A}^T X_\infty + \tilde{C}^T \tilde{C} - X_\infty (B_2 B_2^T - \gamma^{-2} B_1 B_1^T)X_\infty = 0$ has a positive semi-definite solution such that $\tilde{A} - (B_2 B_2^T - \gamma^{-2} B_1 B_1^T)X_\infty$ is stable, where $\tilde{A} = A - B_2 D_{12}^T C_1$ and $\tilde{C}^T \tilde{C} = C_1^T (I - D_{12}D_{12}^T)C_1$.

2. The algebraic Riccati equation $\overline{A}Y_\infty + Y_\infty \overline{A}^T + \overline{BB}^T - Y_\infty (C_2^T C_2 - \gamma^{-2}C_1^T C_1)Y_\infty = 0$ has a positive semi-definite solution such that $\overline{A} - Y_\infty (C_2^T C_2 - \gamma^{-2}C_1^T C_1)$ is stable, where $\overline{A} = A - B_1 D_{12}^T C_2$ and $\overline{BB}^T = B_1(I - D_{21}^T D_{12})B_1^T$.

3. $\rho(X_\infty Y_\infty) < \gamma^2$, where $\rho(\cdot)$ denotes the maximum of the absolute values of the matrix's eigenvalues.

The so-called central controller that solves the suboptimal $H_\infty$ problem can be written in a form that closely resembles the state-estimate feedback form of the LQG controller:

$$\dot{\hat{x}} = A\hat{x} + B_1\hat{w}^\star + B_2 u + \left[ B_1 D_{21}^T + Z_\infty C_{2z}^T \right]\left( y - C_2\hat{x} - D_{21}\hat{w}^\star \right)$$

$$u = -F_\infty \hat{x} \qquad\qquad (100.174)$$

$$\hat{w}^\star = \gamma^{-2} B_1^T X_\infty \hat{x}$$

where $C_{2z} = C_2 + \gamma^{-2}D_{21}B_1^T X_\infty$, $F_\infty = D_{12}^T C_1 + B_2^T X_\infty$, and $Z_\infty = Y_\infty(I - \gamma^{-2}X_\infty Y_\infty)^{-1}$. The dynamic part of the above compensator can be interpreted as an estimate of the state assuming that the worst-case disturbance $w^\star$ is present. The control signal is a linear feedback of the estimated state, just as in the LQG case. Although the controller formulas above look more complicated than in the LQG case, this is largely due to the fact that the controlled variable $z$ has a more general form in the $H_\infty$ problem statement above. It should be noted, however, that unlike the LQG case, the solution of the Riccati equations is coupled in the $H_\infty$ case due to condition 3 above.

## Example

Consider the following state-space system, which represents a plant with two lightly damped modes at approximately $\omega_1 \approx 1$ and $\omega_2 \approx 3.2$:

$$\dot{x}(t) = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -1 & -.1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -10 & -.1 \end{bmatrix} x(t) + \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ .2 & 0 \end{bmatrix} d(t) + \begin{bmatrix} 0 \\ 1 \\ 0 \\ -1 \end{bmatrix} u(t) \qquad (100.175)$$

$$y(t) = \begin{bmatrix} 1 & 0 & .5 & 0 \end{bmatrix} x(t) + \begin{bmatrix} 0 & .01 \end{bmatrix} d(t)$$

Here, the term $d(t) = [d_1(t) \quad d_2(t)]^T$ is a vector whose first term represents the plant disturbance, and whose second term represents the measurement disturbance (deterministic). To pose the LQG control problem, we can propose minimizing the cost function

$$\int_0^\infty \left( x^T T^T T x + r u^2 \right) dt = \int_0^\infty z^T z \ dt \qquad (100.176)$$

where

$$z \triangleq \begin{bmatrix} T \\ 0 \end{bmatrix} x + \begin{bmatrix} 0 \\ r \end{bmatrix} u \qquad (100.177)$$

and $T = [.5 \quad 0 \quad -1 \quad 0]$. However, Eq. (100.175) and (100.177) are also in the form of Eq. (100.172), facilitating an $H_\infty$ design that minimizes the $H_\infty$ norm of the transfer function from the disturbance $d$ to the controlled variable $z$. We can compare this design to an $H_2$ controller design that minimizes the $H_2$ norm (Eq. (100.176)). The two controllers will therefore minimize different norms of the same transfer function.

The results of this comparison are shown in the curves of Fig. 100.55. The distinguishing feature of this comparison is the flattening effect of the $H_\infty$ controller. Although this is a plot of a frequency response magnitude and not the maximum **singular value,** it is apparent that the$H_\infty$ controller is reducing the peak response, while the $H_2$ controller is reducing the average response and providing faster roll-off in the frequency domain.

## Other Approaches

Although the LQG and $H_\infty$ design methodologies are probably the most commonly used techniques for linear MIMO controller design, there are many other optimization-based techniques available. A well-developed theory exists for $L_1$ control, which minimizes the maximum magnification of the peak of the input signal using a linear programming algorithm [Dahleh and Diaz-Bobillo, 1995]. This approach departs from traditional controller design methodologies that have attempted to arrive at a set of formulas for the optimal controller. But with the advent of powerful computers, it makes sense to consider a problem solved if it can be reduced to a problem that can be efficiently solved using a computer algorithm. This is also the premise underlying the design methodology of *linear matrix inequalities*, in which the control design problem is reduced to a convex programming problem [Boyd et al., 1994].
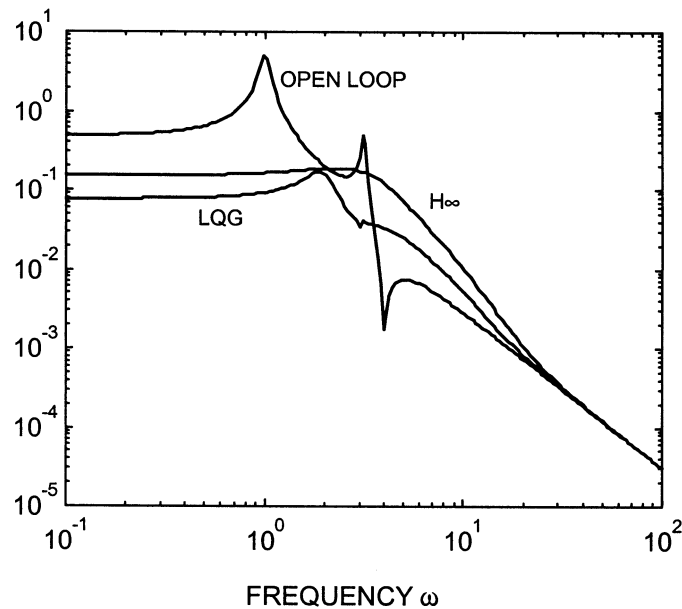
**FIGURE 100.55** Frequency response of the closed-loop transfer function from $d_1$ to $Tx$, comparing the $H_2$ (LQG) and $H_\infty$ control designs, using $r = 0.1$.

## Defining Terms

**Controllability:** A linear system is said to be controllable if a control input exists that will drive a system with an arbitrary initial condition to a desired final state in a finite time.

**Stabilizability:** A linear system is said to be stabilizable if its unstable part is controllable.

**Observability:** A linear system is said to be observable if its state vector can be reconstructed from finite-length observations of its input and output.

**Detectability:** A linear system is said to be detectable if its unstable part is observable.

**Positive-(semi)definite:** A positive- (semi)definite matrix is a symmetric matrix $A$ such that for any nonzero vector $x$, the quadratic form $x^T A x$ is positive (non-negative).

**Observer** (or **estimator**): A linear system whose state output approximates the state vector of a different system, rejecting noise and disturbances in the process.

**Singular value:** Singular values are non-negative real numbers that measure the magnification effect of an operator in the different basis directions of the operator's space.

## References

B. D. O. Anderson and J. B. Moore, *Optimal Filtering,* Prentice-Hall, 1979.

J. S. Bay, *Fundamentals of Linear State Space Systems,* WCB/McGraw-Hill, 1999.

S. Boyd et al., *Linear Matrix Inequalities in System and Control Theory,* Society for Industrial and Applied Mathematics, 1994.

A. E. Bryson, Jr. and Y. C. Ho, *Applied Optimal Control,* Hemisphere Publishing, 1975.

M. Dahleh and I. J. Diaz-Bobillo, *Control of Uncertain Systems: A Linear Programming Approach,* Prentice-Hall, 1995.

J. C. Doyle, K. Glover, P. P. Khargonekar, and B. A. Francis, State-space solutions to standard $H_2$ and $H_\infty$ control problems, *IEEE Trans. on Automatic Control,* AC-34, 831–847, 1988.

B. A. Francis, *A Course in H∞ Control Theory*, Lecture Notes in Control and Information Sciences, Vol. 88, Springer-Verlag, 1987.

M. Green and D. J. N. Limebeer, *Linear Robust Control*, Prentice-Hall, 1995.

R. E. Kalman, A new approach to linear filtering and prediction problems, *Transactions of the ASME, Journal of Basic Engineering*, 82, 35–45, 1960.

D. E. Kirk, *Optimal Control Theory*, Prentice-Hall, 1970.

K. Zhou and J. C. Doyle, *Essentials of Robust Control*, Prentice-Hall, 1998.

## Further Information

Optimal control and estimation is an actively growing field of control systems research. Classic texts in the area include [Kirk, 1970] and [Bryson and Ho, 1975] for optimal control systems, and [Anderson and Moore, 1979] for Kalman filtering, with the original source being [Kalman, 1960]. Further information on $H_2$ and $H_\infty$ theory can be found in [Doyle, et al., 1989], [Zhou and Doyle, 1998], [Green and Limebeer, 1995], and [Francis, 1987].

# 100.9   Neural Control

*Mo-Yuen Chow*

Artificial intelligence had strong ties with automatic control during its early development stages several decades ago. Typical examples of these ties are the development of cybernetics, robotics, and early learning systems. Recent efforts to incorporate aspects of artificial intelligence into the design and operation of automatic control systems have focused on using techniques such as artificial neural networks, fuzzy logic, and expert systems. The application of one or more of these techniques in the design of control systems has come to be known as *intelligent control* [Antsaklis et al., 1994], a term questioned by some for its implications. Whether or not such systems should be classified as intelligent, they represent significant contributions to the field of automatic control, as evidenced by the rapidly growing wealth of literature devoted to the successful application of such systems to complex control problems [Chow and Menozzi, 1994a; 1994b; Chow and Teeter, 1997; Chow and Yee, 1991; Hunt et al., 1992; Miller et al., 1990; Nguyen and Widrow, 1990; Psaltis et al., 1988; Werbos, 1990].

The nonlinear functional mapping properties of Artificial Neural Networks (ANNs) are central to their use in system identification and control [Narendra and Parthasarathy, 1990]. Although a number of key theoretical problems remain, results pertaining to the approximation capabilities of neural networks demonstrate that they have great promise in the modeling and control of nonlinear systems. The artificial neural network technology has become increasingly popular as a tool for performing tasks such as automatic control, system identification, pattern recognition, and time series prediction. Most of the *conventional methods*, such as PI control, are based on mathematical and statistical procedures for the modeling of the system and the estimation of the optimal controller parameters. In practice, the plant to be controlled is often highly nonlinear and a mathematical model may be difficult to derive. In such cases, conventional techniques may prove to be suboptimal and may lack robustness in the face of modeling error, because they are only as accurate as the model that was used to design them. With the advancement of technology, however, sophisticated control using artificial neural network techniques has been developed and used successfully to improve the control of systems that cannot be easily handled by conventional control, thus giving rise to terminology such as *neural control* and *intelligent control*.

Usually, a human operator is responsible for adjusting the controller's parameters in order to use his/her own idea of good performance. Indirectly, the operator is performing a minimization of a cost function based on his/her knowledge. More generally, when confronted with a problem situation, humans execute a mapping between a set of events and the set of corresponding appropriate actions. The appropriateness of these actions is due to some basic acquired knowledge, or even instinct, that guides the initial stages of the mapping. Then, through experience and a set of implicit guidelines, the human learns to perform a better mapping. This is an ongoing process throughout a person's life. Similarly, an ANN, if given initial guidance, can learn to improve its performance through a set of guidelines, (e.g., minimize a cost function). In fact, a properly structured ANN can learn any arbitrarily complicated mapping [Cybenko, 1989; Rumelhart and McClelland, 1986; Werbos, 1974].
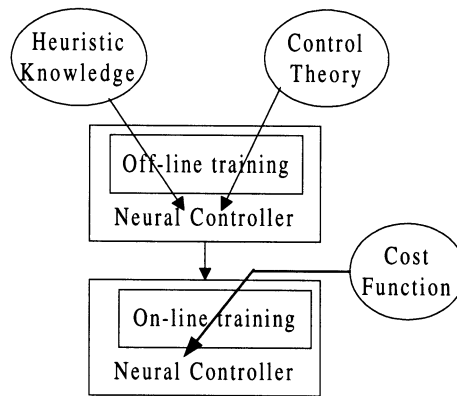
**FIGURE 100.56**    A two-step training process for a neural controller.

Successful adaptation of online neural controllers in many cases requires careful and substantial design effort. One of the common practices is to pre-train a neural controller offline, based on some simplified design methods, in the same way that a PI controller is designed based on *approximated* system models that can provide reasonable control performance, before putting the neural controller for online fine-tuning [Chow and Menozzi, 1993; Chow and Teeter, 1995; Teeter et al., 1996]. This approach, as shown in Fig. 100.56, can speed up the neural controller online adaptation process and increase the closed-loop online adaptation stability because the initial weights of the neural controller are much closer to the optimal final weights (if they exist) after the pre-training process. By learning online, the ANN controller can adapt to changing operating environments.

This chapter section briefly describes the feedforward net paradigm to facilitate the discussion of the *neural observer* and *neural controller* concepts in later sections. An example of using neural control for an HVAC system will then be provided to demonstrate its effectiveness for solving real-world problems.

## Brief Introduction to Artificial Neural Networks

Increasing interest in studying the mechanisms and structure of the brain has led to the development of new computational models for solving problems such as pattern recognition, fast information processing, and adaptation. In the 1940s, McCulloch and Pitts studied the potential and capabilities of the interconnection of components based on a model of a biological neuron. Since then, many different models and architectures have been developed and analyzed for a variety of applications [Zurada, 1992]. One of the most common neuron models is shown in Fig. 100.57.

The inputs *x* to the neuron model are scaled by connection weights *w* and summed. An additional input *b*, often referred to as a *bias*, is added to this sum and the result becomes the input to a function $f(\cdot)$, called the *activation function*, which computes the output of the neuron. The bias can be considered as a connection weight for a constant input of +1. The terms *neuron model* and *neuron* are used interchangeably in this chapter section.

The individual neurons are not very powerful in terms of computation or representation, but the interconnection of neurons to form an *artificial neural network* (ANN) can provide a means of encoding complex relationships between input and output variables. Of the many ANN architectures that have been proposed, the *multilayer feedforward artificial neural network* (MFANN) shown in Fig. 100.58 is one of the most popular.

Bias terms have been omitted in Fig. 100.58 for simplicity. The layers between the input and output layers of an MFANN are usually referred to as *hidden layers* because their inputs and outputs are not measurable at the inputs or outputs of the network. It has been shown that an MFANN with a single hidden layer can approximate any continuous function to an arbitrary degree of accuracy [Cybenko, 1989; Werbos, 1974]. The process of adjusting ANN connection weights in an effort to obtain a desired input/output mapping is usually referred to as *training*. Training represents an optimization problem for which the solution is a set of weights that minimizes some measure of approximation error. The choices of activation functions, number of neurons, error measures, and optimization methods can significantly affect training results.
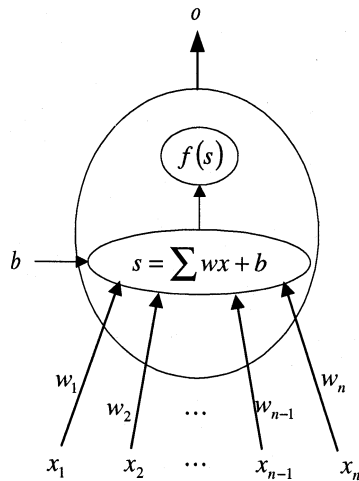
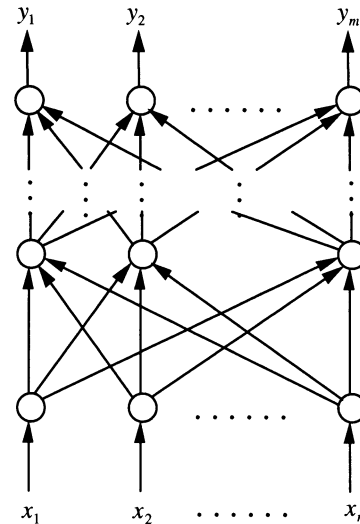**FIGURE 100.57** Computational model of a biological neuron.



**FIGURE 100.58** A multilayer feedforward ANN.

When the desired mapping is described by a set of input/output data, weights are usually modified after the presentation of each input/output pattern. This method is referred to as *pattern update* and represents an approximation of true gradient descent that is generally valid when a sufficiently small stepsize is used. *Batch update*, for which weight changes are accumulated over one sweep of the set of training patterns before being applied, is sometimes used in an effort to more closely mimic true gradient descent. Variants of back-propagation and other training methods can be found in [Zurada, 1992].

## Neural Observer

The nonlinear functional mapping properties of *neural networks* are central to their use in identification and control [Chow and Teeter, 1995; Hunt et al., 1992; Poggio and Girosi, 1990; Teeter and Chow, 1998; Teeter et al. 1994]. Although a number of key theoretical problems remain, results pertaining to the approximation capabilities of neural networks demonstrate that they have great promise in the modeling of nonlinear systems. An important question in system identification is whether a system under study can be adequately represented within a given model structure [Hunt et al., 1992]. In the absence of such concrete theoretical results for neural networks, it is usually assumed that the system under consideration belongs to the class of systems that the chosen network is able to represent. Two system identification techniques are now introduced: *forward modeling* and *inverse modeling*.

The procedure of training a neural network to represent the forward dynamics of a system is often referred to as the *forward system identification* approach [Hunt et al., 1992]. A schematic diagram of this process is shown in Fig. 100.59.

The neural network is placed in parallel with the system, and the error, *e*, between the system outputs, *y*, and network outputs, *ŷ*, is used to train the network. This represents a classical *supervised learning* problem for which the teacher (i.e., the system) provides target values (i.e., system outputs) directly in the output coordinate system of the learner (i.e., the network model) [Jordan and Rumelhart, 1991].

In an *inverse system identification* approach, a network is trained in an effort to model the inverse of the plant mapping [Hunt et al., 1992]. One of the simplest approaches, known as *direct inverse system identification*, is shown schematically in Fig. 100.60.

A synthetic training signal, *s*, is introduced to the system, and the system output, *y*, is used as the input to the network. The network output is compared to the training signal and this error is used to train the network.
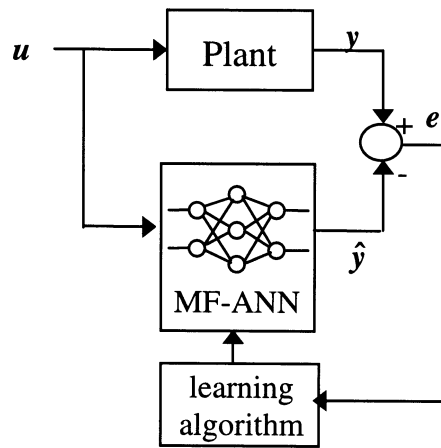
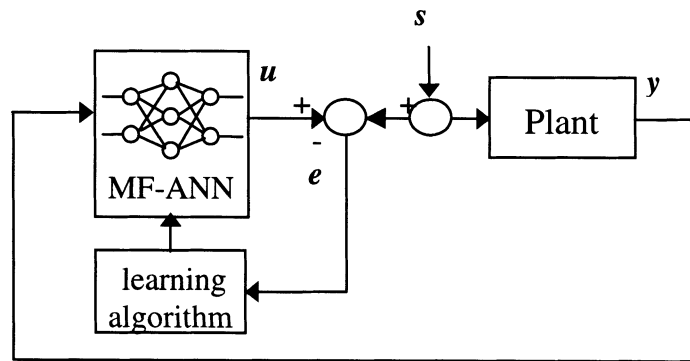**FIGURE 100.59**    Forward system identification approach.



**FIGURE 100.60**    Direct inverse system identification approach.

The inverse modeling structure shown in Fig. 100.60 tends to force the network to represent the inverse of the plant, but there are potential drawbacks to this approach. The training signal must be chosen to sample over a wide range of system inputs, and the actual operational inputs may be difficult to define *a priori* [Jordan and Rumelhart, 1991]. This point is strongly related to the concept of persistent excitation discussed in adaptive control literature. A second drawback is that an incorrect inverse model can be obtained if the nonlinear system mapping is not one-to-one. An approach called *specialized inverse modeling* has been proposed in an effort to overcome these problems. The details of this approach can be found in [Psaltis et al., 1988]. The neural network identification models can be used in the adaptive control of unknown nonlinear plants.

## Neural Control

A method of *direct adaptive control* is depicted in Fig. 100.61.

Methods for directly adjusting control parameters based on the output error $e_c$ are generally not available. This is because the unknown nonlinear plant in Fig. 100.61 lies between the controller and the output error. Until such methods are developed, adaptive control of nonlinear plants must be performed using *indirect* methods [Narendra and Parthasarathy, 1990]. Figure 100.62 depicts a general method of indirect adaptive control using artificial neural networks.

Tapped delay lines (TDL) provide delayed inputs and outputs of the plant to the neural controller and neural observer. Error $e_i$ is used to adapt the neural observer, while the parameters of the neural observer along with
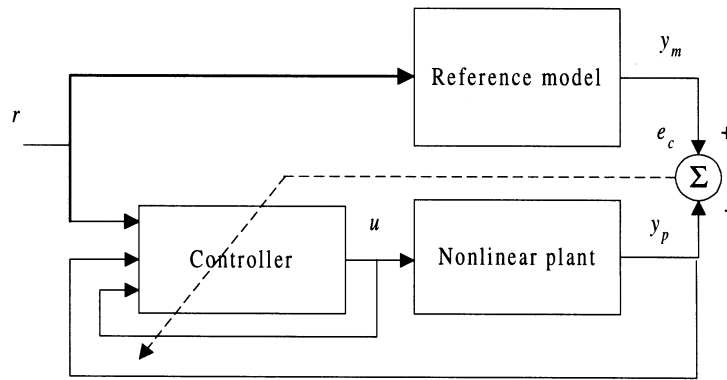
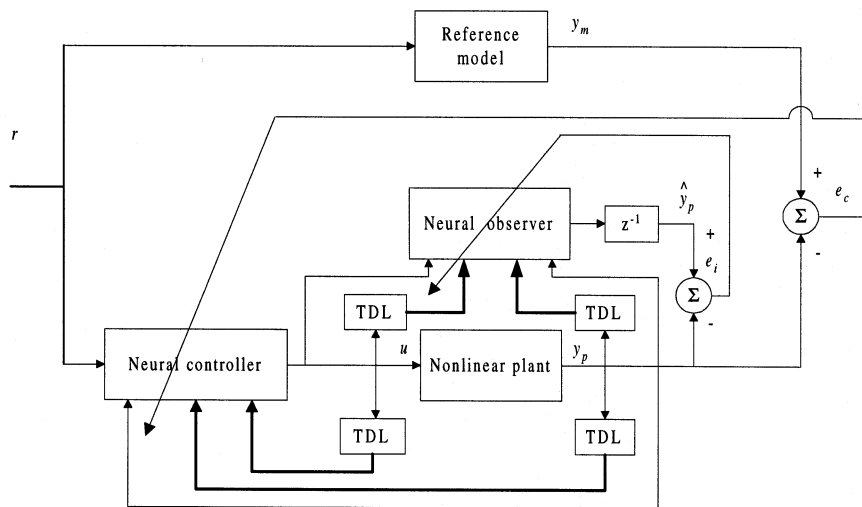**FIGURE 100.61**   Direct adaptive control.



**FIGURE 100.62**   A method of indirect adaptive control using neural networks.
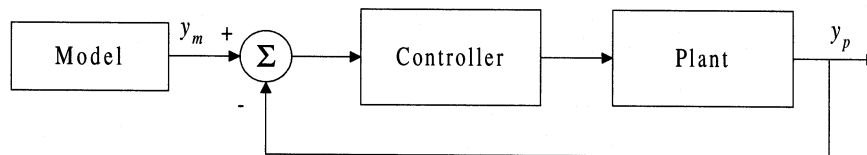


**FIGURE 100.63**   Explicit model-following.

error $e_c$ are used to adapt the neural controller. The model reference approach depicted in Fig. 100.62 is commonly referred to as *implicit model-following* in adaptive control literature [Narendra and Annaswamy, 1989]. This type of model-following is performed when the dynamics of the closed-loop system are forced to asymptotically match the dynamics of the reference model. The method of *explicit model-following* is depicted in Fig. 100.63. In this case, the reference model acts as a prefilter, and the value of the system output is forced to asymptotically track the value of the reference model output [Narendra and Annaswamy, 1989].

A neural control and identification scheme can be described as the following: at each time step, the plant states are measured and a neural network controller computes the plant inputs. Controller parameters are then

adjusted between samples so that the system approaches optimal performance with respect to a given cost index. The general form of the cost index used to adapt the neural controller is:

$$J_k = \sum_{i=k+1}^{N} L\left(\mathbf{x}(i), \mathbf{u}(i)\right) \tag{100.178}$$

where $L$ is the cost function as a function of system state $\mathbf{x}$ and control $\mathbf{u}$. Thus, at each time step, the goal is to minimize $J_k$ subject to system dynamics and control constraints, with $k$ denoting the current time step and $N$ the prediction horizon.

The back-propagation training algorithm can be used to adapt neural networks for the identification and control of nonlinear plants [Teeter and Chow, 1998; Werbos, 1990]. For system identification, a network can be trained offline using plant input/output data obtained from simulation of a mathematical model or from observation of the physical system. When the network is used for adaptive identification, training can be performed using *batch update* with a window of sampled data, or with the *pattern update* method in which training patterns consist of past inputs and outputs measured at each sample time.

In order to adaptively train a neural *controller* using gradient descent, the partial derivatives of a cost index, $J_k$, with respect to the network weights, $\mathbf{w}$, must be obtained [Chow and Yee, 1991]. Let $J_k$ have the form $J_k = L(\mathbf{y}(k+1),\mathbf{u}(k+1)) + L(\mathbf{y}(k+2),\mathbf{u}(k+2)) + \ldots + L(\mathbf{y}(k+n),\mathbf{u}(k+n))$ where $k$ is the current sample. For simplicity of notation, let $L(\mathbf{y}(k+i),\mathbf{u}(k+i))$ be denoted by $L(k+i)$. Application of the chain rule yields

$$\frac{\partial L(k+i)}{\partial \mathbf{w}} = \frac{\partial L(k+i)}{\partial \mathbf{u}(k)}\frac{\partial \mathbf{u}(k)}{\partial \mathbf{w}} = \left[\frac{\partial L(k+i)}{\partial \mathbf{y}(k+i)}\frac{\partial \mathbf{y}(k+i)}{\partial \mathbf{u}(k)} + \frac{\partial L(k+i)}{\partial \mathbf{u}(k+i)}\frac{\partial \mathbf{u}(k+i)}{\partial \mathbf{u}(k)}\right]\frac{\partial \mathbf{u}(k)}{\partial \mathbf{w}} \tag{100.179}$$

The $\partial \mathbf{u}(k)/\partial \mathbf{w}$ term can be calculated using the backpropagation approach since the controller is a neural network. The $\partial \mathbf{y}(k+i)/\partial \mathbf{u}(k)$ and $\partial \mathbf{u}(k+i)/\partial \mathbf{u}(k)$ terms are obtained using the neural controller and identifiers. First, future inputs and outputs of the plant are predicted. The partial derivatives are then obtained by recursively computing the input/output sensitivities of the plant and controller through $i$ samples. This approach is often referred to as *back-propagation through time* [Chow and Yee, 1991; Werbos, 1990].

The training algorithm resembles methods used by Nguyen and Widrow [1990] and others for training a neural controller to achieve an end goal. In this case, however, the output *trajectory* is of interest and the training is performed in realtime (i.e., output values must be repeatedly predicted rather than observed over several trials). A flowchart of the control methodology is shown in Fig. 100.64.

After controller outputs are computed, the weights of the controller are adjusted $N$ times before the next sample time. The value of $N$ can be selected based on time constraints or convergence properties of the neural controller and observers. If $N$ is large, the neural observers are inaccurate; and if a large prediction horizon is used, the adaptation of controller parameters may cause performance to deteriorate.

## HVAC Illustration

In order to demonstrate the ability of the neural identification and control schemes to handle disturbances and changes in nonlinear plant dynamics, a Heat, Ventilation, and Air-Conditioning (HVAC) system where a thermal load is added and the actual output vs. the commanded output of each actuator is modified. The neural observers are adapted at each time step in one simulation, while adaptation of the observers is stopped at time step $k = 200$ in another simulation. Both simulations are performed for 1000 time steps. The reference trajectory is plotted in Fig. 100.65, along with the output of the system that uses nonadaptive identifiers after time step $k = 200$. Tracking errors for both simulations are plotted in Fig. 100.66, where $T_3$ is the temperature to be controlled and $r$ is the reference signal to be tracked.

The performance of the system with nonadaptive observers deteriorates due to the disturbance and the changes in plant dynamics. In this case, adapting the neural observers enables them to more accurately predict
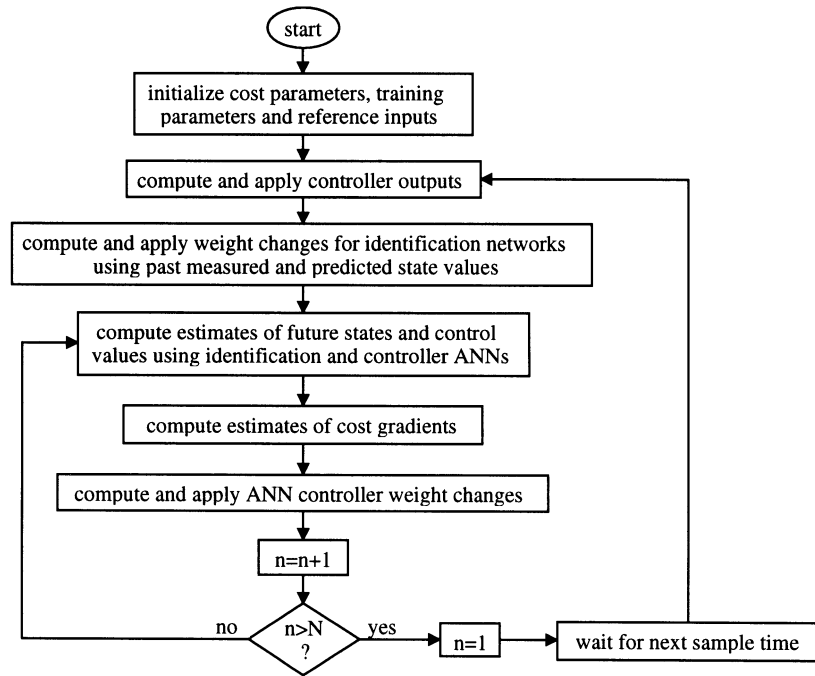
**FIGURE 100.64**    Flowchart of the neural control training methodology.
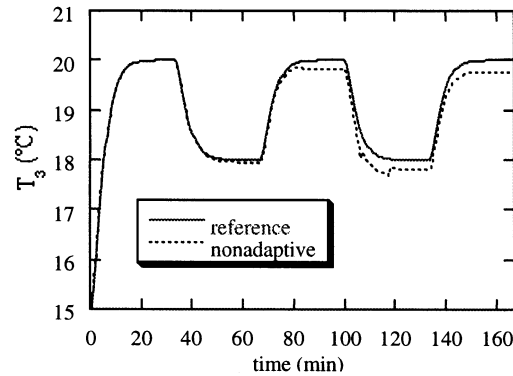


**FIGURE 100.65**    Reference and output trajectories using nonadaptive neural identifiers.

future states and inputs, and results in better tracking. It is important to select appropriate adaptation stepsizes for the observers because only one training pattern, consisting of the most recent set of plant inputs and states, is available for training. In order to compare the neural identification and control methodology with another potential control methodology, a PI-type controller has been designed for the HVAC system. Typical responses of the system with PI-type controller are shown in Fig. 100.67.

The simple PI-type controller satisfies the conservative performance specifications for the cases tested, but does not always use its resources efficiently. For example, if the outside air temperature of the HVAC system is close to the steady-state reference temperature, it may be more efficient to increase the room volumetric flow rate for a period of time in order to reduce the amount of heating or cooling performed by the heat exchanger. The neural and PI-type control schemes are tested using initial conditions $T_3(0) = 15°C$, a constant outside temperature of 22°C, and a steady-state reference value of 20°C. The tracking errors and heat exchanger outputs for both methods are shown in Figs. 100.68 and 100.69, respectively.
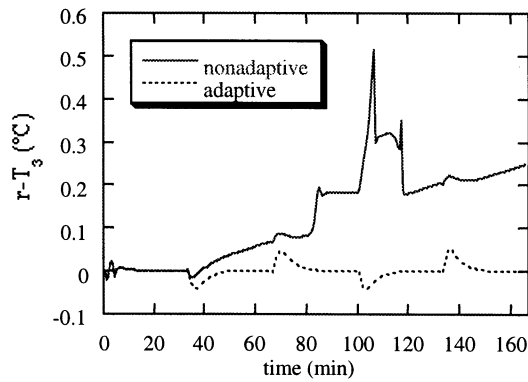
**FIGURE 100.66** Tracking errors for the system using adaptive and nonadaptive neural identifiers.
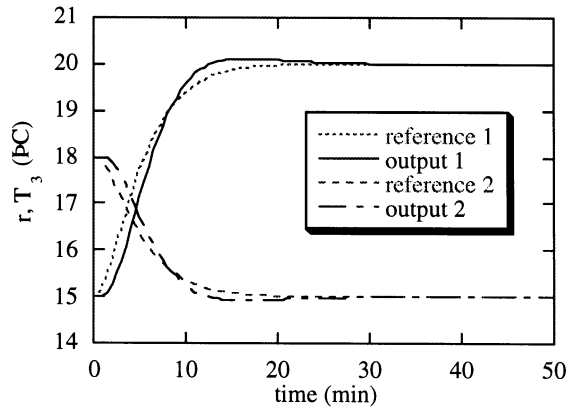


**FIGURE 100.67** Typical responses of the system with PI-type controller.
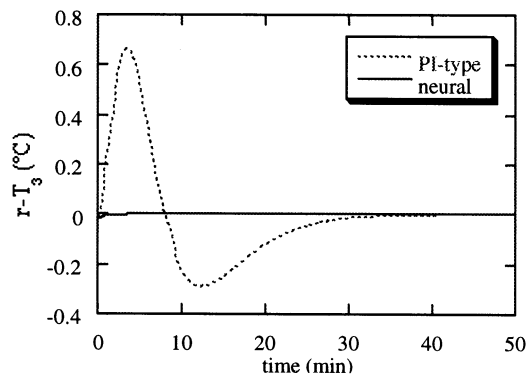


**FIGURE 100.68** Tracking errors of the PI-type and neural control systems.
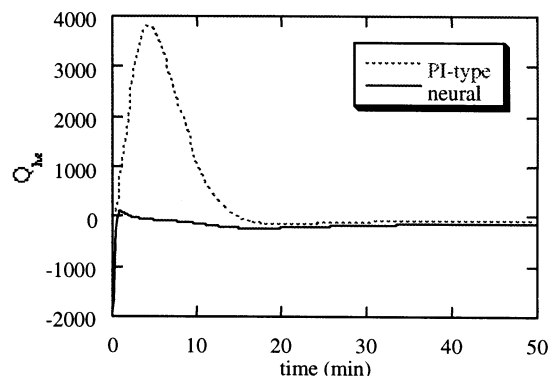


**FIGURE 100.69** Heat exchanger outputs for the PI-type and neural control systems.

## Conclusion

The use of neural networks for system identification and control provides a means of adapting a controller on-line in an effort to minimize a given cost index. The cost index includes typical measures associated with system performance and can be modified without significantly increasing the computational complexity of the adaptation process. For nonlinear systems, the identification networks demonstrate the capacity to learn changes in the plant dynamics. The performance of the neural control and identification methodology compares favorably with many types of conventional approaches.

## References

Antsaklis, P. J., Albus, S., Lemmon, M. D., Mystel, A., Passino, K. M., Saridis, G. N., and Werbos, P. (1994). Defining intelligent control, *IEEE Control Systems,* 4, 5, 58–66.

Chow, M.-Y. and Menozzi, A. (1993). Design Methodology of an Intelligent Controller Using Artificial Neural Networks, *IECON'93,* Maui, Hawaii.

Chow, M.-Y. and Menozzi, A. (1994). A Self-Organized CMAC Controller for Robot Arm Movements, *IEEE International Conference on Industrial Technology,* Guangzhou, China.

Chow, M.-Y. and Menozzi, A. (1994). Using a Cerebellar Model for FES Control of the Upper Limb, *16th Annual International IEEE Engineering in Medicine and Biology Society Conference,* Baltimore, MD.

Chow, M.-Y. and Teeter, J. (1995). A knowledge-based approach for improved neural network control of a servomotor system with nonlinear friction characteristics, *Mechatronics,* 5(8), 949–962.

Chow, M.-Y. and Teeter, J. (1997). Reduced-Order Functional Link Neural Network for HVAC Thermal System Identification and Modeling, *1997 International Conference on Neural Networks,* Houston, TX.

Chow, M.-Y. and Yee, S. O. (1991). An adaptive backpropagation through time training algorithm for a neural controller, *1991 IEEE International Symposium on Intelligent Control,* Arlington, VA, 170–175.

Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function, *Mathematics of Control, Signals, and Systems,* 2, 303–314.

Hunt, K. J., Sbarbaro, D., Zbikowski, R., and Gawthrop, P. J. (1992). Neural networks for control systems — a survey, *Automatica,* 28(6), 1083–1112.

Jordan, M. I. and Rumelhart, D. E. (1991). Forward models: supervised learning with a distal teacher. *Occasional Paper No. 40,* Center for Cognitive Science, MIT.

Miller, III, W. Thomas, Sutton, Richard S., Werbos, Paul J. (1990). *Neural Networks for Control,* The MIT Press, Cambridge, MA.

Narendra, K. S. and Annaswamy, A. M. (1989). *Stable Adaptive Systems,* Prentice-Hall, Englewood Cliffs, NJ.

Narendra, K. S. and Parthasarathy, K. (1990). Identification and control of dynamical systems using neural networks, *IEEE Transactions on Neural Networks,* 1(1), 4–27.

Nguyen, D. and Widrow, B. (1990). The truck backer-upper: an example of self-learning in neural networks, *Neural Networks for Control,* The MIT Press, Cambridge, MA, 287–299.

Poggio, T. and Girosi, F. (1990). Networks for approximation and learning, *Proceedings of the IEEE,* 78(9), 1481–1497.

Psaltis, D., Sideris, A., and Yamamura, A. A. (1988). A multilayered neural network controller, *IEEE Control Systems Magazine,* 17–21.

Rumelhart, D. E. and McClelland, J. L. (1986). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition,* The MIT Press, Cambridge, MA.

Teeter, J. and Chow, M.-Y. (1998). Application of functional link neural network to HVAC thermal dynamic system identification, *IEEE Transactions on Industrial Electronics,* 45(1), 170–176.

Teeter, J., Chow, M.-Y., and Brickley, J. J. Jr. Use of a Fuzzy Gain Tuner for Improved Control of a DC Motor System with Nonlinearities, *IEEE International Conference on Industrial Technology,* Guangzhou, China.

Teeter, J. T., Chow, M.-Y., and Brickley, J. J. Jr.(1996). A novel fuzzy friction compensation approach to improve the performance of a dc motor control system, *IEEE Transactions on Industrial Electronics,* 43(1), 113–120.

Werbos, P. J. (1974). *Beyond Regression: New Tools for Prediction and Analysis in Behavioral Science,* Harvard University Press, Cambridge, MA.

Werbos, P. J. (1990). Backpropagation through time: What it does and how to do it, *Proceedings of the IEEE,* 78(10), 1550–1560.

Zurada, J. M. (1992). *Introduction to Artificial Neural Systems,* West Publishing Company, St. Paul, MN.